

Generalist Foundation Models Are Not Clinical Enough for Hospital Operations

Lavender Y. Jiang^{1,2,3*†}, Angelica Chen^{1†}, Xu Han², Xujin Chris Liu^{2,4},
Radhika Dua^{1,2,3}, Kevin Eaton^{5,6}, Frederick Wolff^{2,7}, Robert Steele^{2,8},
Jeff Zhang^{9,10}, Anton Alyakin^{2,11}, Qingkai Pan², Yanbing Chen^{2,12},
Karl L. Sangwon^{2,5}, Daniel A. Alber^{2,5}, Jaden Stryker²,
Jin Vivian Lee^{2,3,11}, Yindalon Aphinyanaphongs^{6,9,10},
Kyunghyun Cho^{1,3,13}, Eric Karl Oermann^{1,2,3,5,9,14*}

¹Courant Institute School of Mathematics, Computing, and Data Science, New
York University, 60 5th Ave, New York, 10001, NY, USA.

²Department of Neurosurgery, NYU Langone Health, 550 First Avenue, New
York, 10016, NY, USA.

³Global AI Frontier Lab, New York University, 1 Metrotech Center, Fl. 22,
Brooklyn, 11201, NY, USA.

⁴Electrical and Computer Engineering, Tandon School of Engineering, 6
MetroTech Center, Brooklyn, 11201, NY, USA.

⁵Grossman School of Medicine, New York University, 550 First Avenue, New
York, 10016, NY, USA.

⁶Department of Medicine, NYU Langone Health, 550 First Avenue, New York,
10019, NY, USA.

⁷Department of Computer Science, ETH Zurich, Universitätsstrasse 6, City,
8092, Zurich, Switzerland.

⁸Department of Surgery, NYU Langone Health, 1 Park Avenue, New York,
10016, NY, USA.

⁹Division of Applied AI Technologies, NYU Langone Health, 227 East 30th
Street, New York, 10016, NY, USA.

¹⁰Department of Population Health, NYU Langone Health, 450 First Avenue,
New York, 10019, NY, USA.

¹¹School of Medicine, Washington University of St. Louis, 660 S. Euclid Ave.,
St. Louis, 63110, MO, USA.

¹²School of Global Public Health, New York University, 708 Broadway, New
York, 10003, NY, USA.

¹³Prescient Design, Genentech, 149 5th Ave. 3rd floor, New York, 10019, NY, USA.

¹⁴Department of Radiology, NYU Langone Health, 450 First Avenue, New York City, 10019, NY, USA.

*Corresponding author(s). E-mail(s): Lavender.Jiang@nyu.edu;
Eric.Oermann@nyulangone.org;

Contributing authors: Angelica.Chen@nyu.edu; Xu.Han@nyulangone.org;
xl3942@nyu.edu; rd3571@nyu.edu ; Kevin.Eaton@nyulangone.org;
wfrederick@ethz.ch; Robert.Steele@nyulangone.org;
Jeff.Zhang@nyulangone.org; Anton.Alyakin@nyulangone.org;
ivrinom@gmail.com; yc6785@nyu.edu; Karl.Sangwon@nyulangone.org;
Daniel.Alber@nyulangone.org; Jaden.Stryker@nyulangone.org;
jin.v.lee@wustl.edu; yin.a@nyulangone.org; Kyunghyun.Cho@nyu.edu;

[†]These authors contributed equally to this work.

Abstract

Hospitals and healthcare systems rely on operational decisions that determine patient flow, cost, and quality of care. Despite strong performance on medical knowledge and conversational benchmarks, foundation models trained on general text may lack the specialized knowledge required for these operational decisions. We introduce **LANG1**, a family of models (100M–7B parameters) pretrained on a specialized corpus blending 80 billion clinical tokens from NYU Langone Health’s electronic health records (EHRs) and 627 billion tokens from the internet. To rigorously evaluate LANG1 in real-world settings, we developed the REalistic Medical Evaluation (**REMedE**), a benchmark derived from 668,331 EHR notes that evaluates five critical tasks: 30-day readmission prediction, 30-day mortality prediction, length of stay, comorbidity coding, and predicting insurance claims denial. In zero-shot settings, both general-purpose and specialized models underperform on four of five tasks (36.6%–71.7% AUROC), with mortality prediction being an exception (up to 94.2% AUROC). After finetuning, LANG1-1B **outperforms** finetuned generalist models up to $70\times$ larger and zero-shot models up to $671\times$ larger, improving AUROC by 3.64%–6.75% and 1.66%–23.66% respectively. We also observed cross-task scaling with joint finetuning on multiple tasks leading to across the board improvement on other tasks. LANG1-1B effectively **transfers** to out-of-distribution settings, including other clinical tasks and an external health system. Our findings suggest that predictive capabilities for hospital operations **require** explicit supervised finetuning, and that this finetuning process is made more efficient by in-domain pretraining on EHR. Our findings support the emerging view that specialized LLMs can compete with generalist models in specialized tasks, and show that effective healthcare systems AI requires the combination of in-domain pretraining, supervised finetuning, and real-world evaluation beyond proxy benchmarks.

Keywords: pretraining, finetuning, Electronic Health Records, hospital operations, clinical prediction tasks, domain-specific models

1 Main

Healthcare systems face high-stakes operational decisions daily: which patients are at imminent risk of decline, who can be safely discharged, how many beds will be available for new admissions. These decisions directly impact resource allocation, care coordination, and patient outcomes [1–3]. As healthcare systems face increasing patient volume, there is a need for tools that can analyze complex clinical data to inform these critical decisions. Foundation models, with their powerful text comprehension capabilities and versatility in specialized domains [4–9], have emerged as a promising technology for optimizing hospital operations.

However, deploying LLMs in clinical settings is fraught with challenges. While LLMs show promise in various clinical tasks [10–25], there is disagreement on whether smaller specialized models ("specialists") can outperform general-purpose models ("generalists") [26–28]. Many evaluations rely on proxy benchmarks that weakly reflect real-world clinical constraints like data scarcity and temporal shifts [17, 29–35]. Data privacy concerns further limit the pretraining data for the vast majority of the clinical LLMs [36] to a small set of public corpora, MIMIC [37] and PubMed [38], even though large-scale EHR datasets are known to improve out-of-domain generalization [28, 39–42].

In this work, we focus on hospital operation tasks that represent the daily challenge of healthcare delivery and explore the tradeoffs between off-the-shelf generalists and specialized models trained on a health system’s internal patient notes. Our contributions are as follows:

1. **LANG1: a family of models specialized for hospital operations.** We present LANG1, a suite of decoder LLMs (100M, 1B, and 7B), pretrained from scratch on a mix of 80 billion tokens of EHR notes and 627 billion tokens of internet texts. After task-specific finetuning, LANG1 **outperforms** both off-the-shelf generalist LLMs (such as DEEPSEEK R1 671B) and the parameter-efficient finetuned variant (LoRA finetuned DEEPSEEK DISTILLED LLAMA 70B), on the REMeDE benchmark. Instruction finetuned on one or more tasks, LANG1 **transfers** zero-shot to related tasks and to a different hospital, surpassing generalist models of similar scales.
2. **REMeDE: an operations-grounded evaluation suite.** We construct an internal evaluation benchmark consisting of five clinically important classification tasks from a multi-hospital academic health system across 10 years, where each task contains 87,974 to 421,429 real patients. The benchmark has time-based splits to mimic deployment settings and data-efficiency protocols to reflect real-world constraints.
3. **Engineering Principles for clinical utilities.** We analyze the training dynamics and show that pretraining on next token prediction of unlabeled clinical notes and web text contributes to emergent skills on comprehension tasks but is insufficient for excelling on REMeDE, which specifically requires supervised finetuning (SFT). However, SFT is made more efficient by in-domain pretraining, and larger models pretrained on more clinical data improve temporal robustness.

Overall, our results suggest that health systems with the capacity for in-house model development can gain clear advantages from smaller specialized models, providing a practical and data-efficient pathway to robust operational prediction with minimal task-specific supervision.

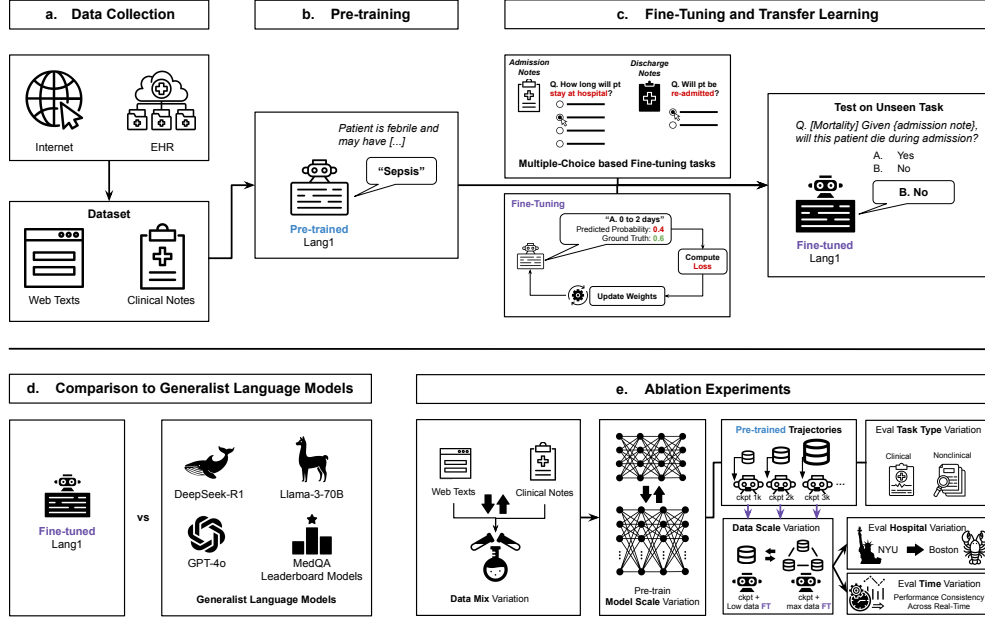


Fig. 1: Overview. (a) We mix unlabeled EHR notes and web texts as our pretrain corpus. (b) We pretrain using next token prediction. (c) Instruction finetuning in multiple choice format enables cross-task transfer. (d) We compare Lang1 to off-the-shelf generalist models. (e) In order to derive design principles, we do ablations on data mix, model scale, pretrain trajectories, data scale, eval task type, eval hospital, and eval time.

1.1 Overview

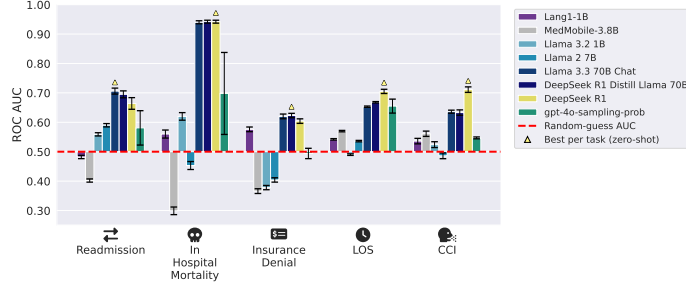
Our work consists of five stages: data collection, pretraining, finetuning, evaluation, and ablations. As shown in Figure 1, we first collect unlabeled clinical notes from NYU Langone EHR and web texts from the internet and mix them to form the corpus (Figure 1a). We pretrain LAng1 via next token prediction (Figure 1b). We instruction finetune LAng1 to predict labels for real-world clinical tasks, enabling cross-task transfer (Figure 1c). We then compare finetuned LAng1 with off-the-shelf generalist models (Figure 1d) and perform ablations of the data mix, model scale, eval task type, pretraining trajectory, finetune data scale, eval data time, and eval hospital to derive design principles for clinical utilities (Figure 1e).

We evaluate LAng1 using ReMeDE, an internal benchmark of real-world, high-impact clinical tasks beyond diagnosis. Unlike recent benchmarks that focus on multi-turn diagnostic dialogue [15, 24, 29], which captures an important but narrow part of clinical decision making, ReMeDE is based on 668,331 EHR notes and emphasizes operational tasks that better represent the day-to-day challenges of healthcare delivery [1–3]. These tasks support practical goals such as reducing costs, optimizing resource use, and improving continuity of care.

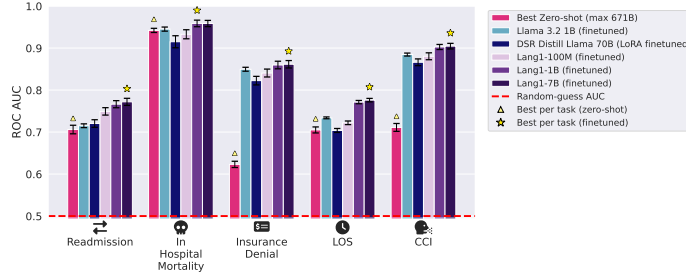
ReMeDE includes five predictive tasks drawn from real-world hospital workflows: 30-day all-cause readmission, in-hospital mortality, binned length of stay, insurance denial, and imputation of binned Charlson Comorbidity Index (CCI, a measure of patient comorbidity burden). These tasks reflect critical decisions tied to patient outcomes, resource planning,

and healthcare operations (See Methods 5.2.2 for more details). To assess model robustness to temporal distribution shifts, each task is evaluated across three non-overlapping test splits drawn from distinct time periods (Appendix A). ReMedE supports flexible few-shot evaluation across a range of language model interfaces and task formats. It is easily extensible for new tasks and evaluation settings. We plan to release ReMedE as a secure evaluation service, allowing trusted researchers to submit models and receive standardized evaluation results without direct access to patient data. This design safeguards patient privacy while enabling fair and reproducible model comparison.

2 Results



(a) Both generalists and specialists underperform zero-shot. Yellow triangles indicate the best zero-shot performance per task. While the best mortality prediction AUROC is 94.2%, performance of other tasks (readmission, insurance denial, LOS, CCI) range from 36.6%–71.7% AUROC.



(b) Finetuned LANG1-1B (purple) outperform best zero-shot performance (magenta) by 1.66% to 23.66% AUROC and finetuned LLAMA 3.2 1B (light blue) and LoRA finetuned LLAMA 70B (deep blue) by 3.64% to 6.75% AUROC. Yellow stars indicate the best finetuned performance per task.

Fig. 2: Finetuned small specialists outperform strong generalists on ReMedE.

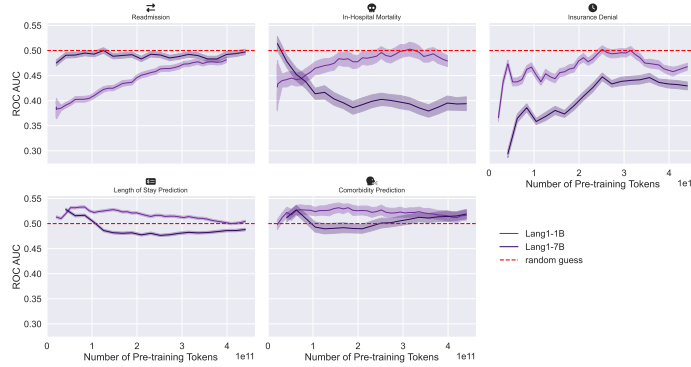
Large generalist models underperform on real-world clinical predictive tasks.. We evaluate both large generalist foundation models and medQA leaderboard models under zero-shot inference, and find that they underperform on ReMedE tasks. After finetuning, LANG1-1B outperform both two finetuned models (LLAMA 3.2 1B and LoRA DEEPSEEK R1 DISTILL LLAMA 3.2 70B) by 3.64% to 6.75% AUROC. LANG1-1B also outperforms the best zero-shot performance (including models up to DEEPSEEK R1 671B) by 1.66% to 23.66% AUROC.

Figure 2a shows that on specialized tasks such as insurance denial, both leaderboard models (LLAMA 2 7B—light blue, LLAMA 3.2 1B—blue, MEDMOBILE—grey) and our own pretrained models (LANG1-1B—purple) underperform in the zero-shot setting. While mortality prediction has up to 94.2% AUROC, the other four tasks range between 36.6%–71.7% AUROC. This shows that even the strongest generalist models falter on these specialized operational tasks.

Figure 2b compares the best zero-shot result per task against a few finetuned models, including LANG1 (100M, 1B and 7B), LLAMA 3.2 1B, and a parameter-efficient finetuned version of DEEPSEEK R1 DISTILL LLAMA 70B. ACROSS all five tasks, LANG1-1B and LANG1-7B (purple bars) achieve higher AUROC than the best zero-shot model (magenta) and the other finetuned models (light blue and dark blue). The improvements over the other finetuned models range from 3.64%–6.75%. The improvements over the best zero-shot baseline range from 1.66% (mortality) to 23.66% (insurance denial), with the largest gains observed at insurance denial prediction, which is DEEPSEEK R1’s worst performing task.



(a) Reading comprehension performance increases from pretraining.

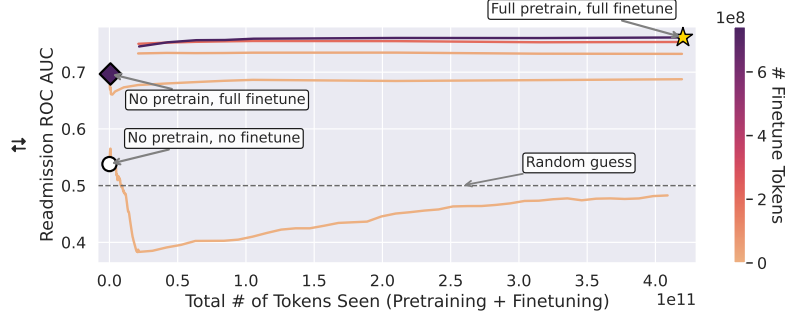


(b) Clinical classification performance does not rapidly emerge from pretraining.

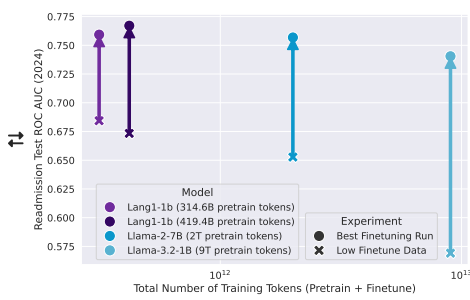
Fig. 3: Zero-shot clinical classification performance does not increase over the course of pretraining, unlike reading comprehension. Error bars depict the 95% confidence interval.

Clinical performance does not emerge during pretraining, unlike reading comprehension. We tracked zero-shot performance of LANG1 (1B and 7B) throughout pretraining, as a function of the number of tokens seen. On comprehension tasks (Method 5.2.4 for data details), accuracy increased with additional pretraining data (Figure 3a), consistent with the intuition that language models improve on text-based reasoning tasks as they are exposed to more data. In contrast, zero-shot AUROC on clinical classification tasks (REMEDe) remained close to or

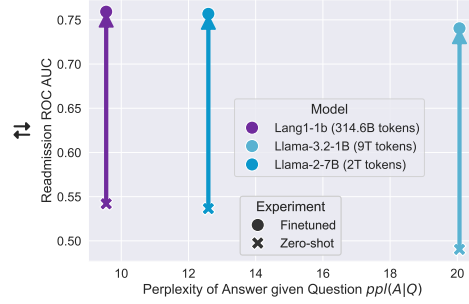
below random chance across the entire pretraining trajectory (Figure 3b). We hypothesize that the mapping from clinical notes to outcomes does not emerge from learning next token prediction on unlabeled texts alone, but must be learned through either task-specific finetuning or alternative pretraining objectives.



(a) Finetuning is more token-efficient for performance gains, but pretraining still provides value. At any fixed token budget (a vertical slice on the x axis), using more finetuning tokens (darker colors) yields a higher ROC AUC for LANG1-1B’s checkpoint trajectory. Nonetheless, a clear gap remains between fully finetuning without pretraining (purple diamond) and the fully pretrained model (yellow star).



(b) Clinically pretrained models (purple LANG1-1B), when finetuned, outperform generalist models of similar size (blue LLAMA models), especially in low data regime (cross). Arrows track the same pretrained model as it is finetuned on different numbers of examples.

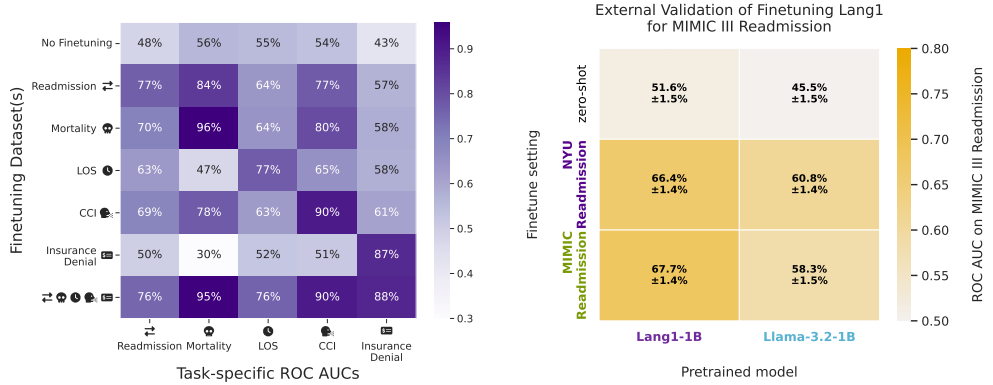


(c) Lower model perplexity on REMeDE answer-question pairs is associated with better zero-shot (cross) and finetuned (circle) performance. LANG1 (purple) has lower perplexity despite fewer total pretrain tokens (though more clinical tokens). Arrows track the same model’s performance.

Fig. 4: Finetuning is are more token-efficient than pretraining for performance gains (Figure 4a), but in-domain pretraining enables sample-efficient finetuning (Figure 4b). This advantage is also associated with lower perplexity on in-domain tasks (Figure 4c).

Compute is better spent on finetuning, but pretraining makes finetuning more efficient. Figure 4a examines the pretraining and finetuning trajectory of the LANG1-1B readmission model under a fixed total token budget (pretraining + finetuning, *i.e.*, each vertical slice). During pretraining of LANG1-1B, checkpoints were saved after each one million training tokens.

Each pretrain checkpoint is finetuned using 100–362,259 discharge notes with readmission label (2.0M–742.0M tokens). The compute budget is the sum of pretraining token for that checkpoint and the number of finetuning tokens used to finetune that particular checkpoint. Within each slice, increasing the proportion of finetuning tokens (darker colors) consistently improves performance. At the same time, **pretraining still provides value**: even with maximal finetuning data, models initialized from pretraining (purple diamond) outperform the randomly initialized one (yellow star) by 7.03% AUROC. A similar pattern appears in Figure 4b: when finetuning data are scarce, LANG1-1B (purple) outperforms generalist models of comparable scale (blue LLAMA-2-7B and LLAMA-3.2-1B) that were pretrained on more nonclinical tokens, demonstrating the efficiency gains of domain-specific pretraining. In fact, Figure 4c shows that LANG1-1B, despite being trained on fewer tokens, achieves lower perplexity on answer–question pairs and stronger zero-shot and finetuned performance than LLAMA-2-7B and LLAMA-3.2-1B. Ablations show that larger models trained on more recent data have better performance (Appendix D).



(a) Finetuning LANG1-1B can often transfer across tasks. The heatmap shows the finetuned model’s performance when finetuned on a subset of ReMedE tasks (y axis) and evaluated on all five ReMedE tasks (x axis).

(b) Finetuning on NYU Readmission (purple) transfers well (darker yellow) to MIMIC III Readmission (green). Overall performance is better on finetuning LANG1-1B (purple) compared to LLAMA 3.2 1B (blue).

Fig. 5: LANG1 is able to transfer to unseen task (Figure 5a) and a different health system (Figure 5b).

LANG1 is able to transfer to unseen task and a different health system. The heatmaps in Figure 5a show LANG1-1B finetuned on one or all ReMedE tasks (rows) and evaluated on all five tasks (columns). Overall, Lang1-1B achieves strong single-task (diagonal) and joint-task (last row) performance, and is well calibrated (Appendix E). Many tasks transfer. For example, finetuning on readmission (second row) boosts performance on the other four tasks. However, this transfer could be asymmetric, *i.e.*, mortality helps LOS (third row, third column), but LOS does not help mortality (fourth row, second column), which can be explained using domain knowledge (Appendix F). Compared to LANG1-1B, LLAMA-3.2-1B has overall worse performance and a different transfer pattern (Appendix C), suggesting that the pattern depends on the pretrained model.

Figure 5b shows LANG1-1B and LLAMA-3.2-1B finetuned using different readmission data (MIMIC III or NYU) and test on MIMIC. Finetuning LANG1-1B shows better performance on both data. For LANG1-1B, finetuning on MIMIC is slightly better than NYU by 1.2% AUROC. However, for LLAMA-3.2-1B, finetuning on NYU is surprisingly better than finetuning on MIMIC by 2.5% AUROC. We suspect this is because NYU+ Readmission has more labeled pairs than MIMIC readmission, suggesting that nonclinical models may benefit more from a large number of slightly out-of-distribution examples compared to a smaller number of in-distribution examples. Indeed, downsampling NYU dataset to the same size as MIMIC would lead to similar results for LLAMA-3.2-1B (Figure M11b). Appendix M extends this analysis to MIMIC mortality and LOS with similar observations.

3 Discussion

We present what is, to our knowledge, the first comprehensive study of LLMs for hospital and health-system operations. We detail a full-cycle approach, in which we build, benchmark and test the robustness of these models as part of the their operational deployment within the NYU Langone Health System. Our work has implications for the broader debate between generalist and specialist models, LLM generalizability, and pretraining-finetuning dynamics.

Why Operational Tasks Matter. Much of the current excitement in medical AI centers on diagnostic reasoning [12, 29, 35, 43–45]. These are valuable directions, but they do not fully capture the day-to-day challenges physicians face. Healthcare is as much operational as it is clinical. Physicians spend only 26% of their time in direct patient care, with much of the remainder devoted to documentation, insurance, and scheduling [1]. Operational outcomes such as readmission, insurance denial, and length of stay directly shape costs, capacity, and continuity of care. Predicting them is actionable, and improving them is measurable. By advancing operational tasks alongside diagnostic ones, we can reduce costs, improve care delivery, and make healthcare more accessible. Notably, we find that many of these tasks require specialized finetuning, and are not out-of-the-box accessible to generalist models, likely reflecting their poor representation in web-scale datasets.

On The Need for Real-World Evaluation. Our results show the limitations of relying on proxy benchmarks as evidence of readiness for clinical deployment. Several MedQA leaderboard models under-perform on ReMedE, suggesting that proxy benchmark success does not establish clinical utility. Evaluating models directly on real-world, task-specific outcomes is therefore essential, not only to identify models that genuinely improve hospital operations, but also to avoid overestimating the safety or value of systems based solely on proxy measures [17, 33, 34].

The Costs of Training. A common concern about specialized models is whether they can be trained affordably outside of industrial labs. Our experience with LANG1 suggests that the answer is yes. Training the 1B-parameter model on 314.5B tokens (150k steps) required roughly 30 days on 64 H100s, which at AWS p5.48xlarge pricing corresponds to a cost of about \$180,000. While non insignificant, this figure is orders of magnitude lower than the multimillion-dollar budgets required for frontier generalist models [5, 6, 46, 47], yet

delivers substantial performance gains on real-world hospital tasks. For a large health system, such costs are comparable to routine IT infrastructure upgrades and therefore financially and operationally feasible and profitable. This supports recent position papers [48, 49] that argue small, specialized models are the future of agentic AI because they can be more reliable, economical, and aligned with domain-specific needs.

In-House Models. From an operational standpoint, training and deploying models in-house can make more sense than relying exclusively on generalist models delivered via commercial APIs. In this paradigm, hospitals not only reduce long-term costs and safeguard patient data, but also retain the ability to adapt models as documentation practices, coding standards, and patient populations shift over time, a need highlighted by our temporal robustness results. In contrast, the dependence on external APIs can introduce ongoing expenses, privacy risks, and limited control over model behavior. Our findings suggest that even modestly sized models, trained locally, can outperform larger off-the-shelf alternatives, reframing clinical AI from “renting intelligence” to “building institutional assets” that evolve with the health system itself.

Pretraining Is Not All You Need. While the reading comprehension capabilities emerge directly from large-scale pretraining, our finding provides evidence that high-stake objective predictions represent a different class of problem. Strong performance on ReMeDE tasks require explicit finetuning and does not emerge from pretraining alone, even with domain-specific data. This reliance on finetuning is shared by general-purpose chatbots, but the underlying tasks are fundamentally different. Chatbot finetuning aligns emergent generalist skills to *subjective*, preference-based goals [50]. In contrast, ReMeDE tasks require finetuning to build a new, non-emergent predictive skill against an *objective*, ground-truth target. We hypothesize that this is because clinical predictions rely on complex relationships not well captured by a next-token-prediction objective. Our experiments show that these tasks demand domain-specific supervision built on top of in-domain pretraining.

Why Transfer Matters. Healthcare tasks often suffer from limited labels due to the expertise required for annotation, and in some cases it is practically impossible to obtain large labeled datasets (*e.g.*, for rare conditions). Our experiments show that instruction finetuning enables transfer across related tasks, where supervision on one outcome (*e.g.*, readmission) improves performance on others (*e.g.*, mortality, length of stay). We also showed that finetuning LANG1 on NYU can transfer to a different hospital. This capability is especially valuable in clinical settings because it reduces dependence on costly annotation pipelines and makes it possible to tackle tasks where labels are scarce or unattainable, effectively maximizing the value of individual annotations.

4 Conclusions

We introduced LANG1, a family of domain-specialized models trained on EHR and web text. To evaluate LANG1, we developed ReMeDE, an operations-grounded benchmark for evaluating language models across five operationally significant tasks. We find that large generalist models underperform LANG1-1B by 1.66% - 23.66% AUROC, despite their strong performance on general and proxy medical benchmarks. LANG1-1B also outperformed other

finetuned models up to $70\times$ larger by 3.64%–6.75%. These results demonstrate that success on proxy benchmarks does not guarantee effectiveness in deployment-critical settings.

Our analysis shows that clinical prediction capabilities do not arise from pretraining alone; explicit finetuning is necessary, though in-domain pretraining improves data efficiency. LANG1 also demonstrates promising transfer across related tasks and health systems, suggesting that carefully designed instruction finetuning can reduce dependency on expensive annotation. Training costs for LANG1 are substantially smaller than frontier models and affordable for large health systems, making specialized models both effective and practical.

Future works include expanding REMEDe to additional institutions, more diverse patient populations, and task types beyond classification outcomes. This is important for establishing shared standards for evaluating models in real operational contexts.

Our findings challenge the assumption that ever-larger internet-trained models will generalize to all domains and instead point to a more hopeful direction for clinical AI. In healthcare, where patient safety is paramount, progress must be deliberate and grounded in real operational outcomes. Benchmarks like REMEDe and specialized models such as LANG1 show that smaller, domain-specific systems can be accurate, affordable, and adaptable, offering a scalable path forward that directly supports hospital operations and improves patient care.

5 Methods

5.1 Data Collection

Data are extracted via structured query language (SQL) scripts to query the NYU Langone EHR, prototyped in an interactive web-based editor (Cloudera Hue), and exported as comma-separated files (CSVs) to an on-prem high-performance computing (HPC) cluster.

Preprocessing. Raw CSV notes (including pathology, radiology, and general hospital notes) are loaded with standard ASCII-encoding using Python Dask [51] for distributed processing. We concatenate narrative fields, standardize punctuation, spacing and formatting via regular expression substitutions, remove non-ASCII and malformed characters, remove errant whitespace and newlines, and filter out short notes (less than 10 words or with placeholder values such as <NA>).

5.2 Datasets

5.2.1 Pretraining Dataset

Web texts. We use SlimPajama (627B tokens) [52], a large extensively deduplicated, multi-corpora, open-source dataset for training LLMs. Its sources include Commoncrawl [53], C4 [54], Github, Books [55, 56], Arxiv, Wikipedia and StackExchange.

NYU Notes. This dataset is previously created using unlabeled inpatient hospital notes signed by medical professionals from the NYU Langone Health EHR¹ for patient encounters starting from January 2011 to May 2020. NYU Notes contains 387,144 patients, 7,247,694 notes, and 4,112,249,482 words in total. NYU Notes are used to train and evaluate NYUTron. [42]

NYU Notes+. This dataset builds on NYU Notes by including a wider range of note types and covering a longer time span, resulting in a total word count 14.5 times greater than that of NYU Notes. NYU Notes+ contains unlabeled hospital notes, pathology note and radiology notes from NYU Langone Health EHR from 2003 to 2023. It comprises 11,689,342 patients, 180,487,092 notes, and 59,917,646,788 words.

5.2.2 Finetuning Datasets and ReMedE Test Set

We derive five task-specific labeled datasets by combining NYUTron [42] finetune datasets, with the addition of 2024 temporal test set to approximate deployment robustness. The 2024 temporal tests set are used for ReMedE. See [Appendix A](#) for a visualization of the data split timeline and [Appendix B](#) for detailed dataset statistics. [Appendix L](#) shows that a small percentage of patient overlap does not over-estimate model performance on readmission. For both zero-shot evaluation and finetuning (section 5.4), the dataset were converted to multiple choice format ([Appendix H](#)).

¹This study is approved by the Institutional Review Board (IRB) at NYU Langone Health. The methods are carried out in accordance with the IRB’s relevant guidelines and regulations.

NYU+ Readmission. Readmission occurs when a patient returns to the hospital shortly after discharge. Predicting readmissions is critical for identifying patients who need longer stay or post-discharge support, and it also serves as a key hospital quality metric. This finetuning dataset contains discharge notes with 30-day all-cause readmission labels. The notes comprise of a subset of NYU+ Notes whose encounters end between January 2013 and November 2021, and additional discharge notes from 2024 for the temporal test. Rehabilitation, dialysis and palliative care notes are excluded to focus on modeling acute readmission. A positive label is assigned if the patient is readmitted within 30 days of discharge, and a negative label otherwise. We split the dataset into five sets. The first three sets are train, validation, and test set with a 8:1:1 ratio from 2013 to May 2021. The 2021 temporal test includes notes from June to December 2021. The 2024 temporal test set includes notes from 2024. The positive class ratio range from 10.81% to 11.29% across these five sets. The dataset contains 421,429 patients, 604,326 notes and 607,877,177 words in total.

NYU+ In-Hospital Mortality. In-hospital mortality prediction identifies patients at highest risk of death during their admission, enabling timely palliative care consultations and goals-of-care discussions that align treatment with patient prognosis. This finetuning dataset contains History and physical (H&P) notes with in-hospital mortality labels. The notes comprise of a subset of NYU+ Notes for encounters ending between January 2013 and November 2021, with additional H&P notes from 2024 for the temporal test. A positive label is assigned if the discharge disposition is "Expired", and a negative label otherwise. We split the dataset into five sets. The first three sets are train, validation, and test set with a 8:1:1 ratio from 2013 to May 2021. The 2021 temporal test includes notes from June to December 2021. The 2024 temporal test set includes notes from 2024. The positive class ratio range from 1.78% to 1.93% across these five sets. In total, the dataset contains 395,991 patients, 566,748 notes, and 608,603,182 words.

NYU+ Length of Stay (LOS). LOS is the number of days a patient remains hospitalized. Predicting LOS is essential for bed management, staffing allocation, and discharge planning. This finetuning dataset contains H&P notes with label for binned length of stay. The dataset comprises of a subset of NYU+ Notes for encounters ending between January 2013 and November 2021, with additional H&P notes from 2024 for the temporal test. We assign the labels based on quantile. We assigned label 0 for an LOS below the 25% quantile (0 to 2 days), label 1 for an LOS between the 25% and 50% quantile (3 days), label 2 for an LOS between the 50% and 75% quantile (4 to 5 days) and label 3 for an LOS above the 75% quantile (> 5 days). We split the dataset into five sets. The first three sets were train, validation, and test set with a 8:1:1 ratio from 2013 to May 2021. The 2021 temporal test includes notes from June to December 2021. The 2024 temporal test set includes notes from 2024. The majority class ratio (0 to 2 days) range from 41.49% to 45.64% across these five sets. The minority class ratio (more than 5 days) range from 23.92% to 26.34%. In total, the dataset contains 395,991 patients, 566,748 notes and 608,603,182 words.

NYU+ Insurance Denial. Insurance denials occur when payers reject claims for hospital services. Predicting denials allows hospitals to proactively address documentation gaps, reducing administrative burden and preventing unexpected out-of-pocket costs for patients.

This finetuning dataset contains H&P notes with insurance denial label. The notes comprise of a subset of NYU+ Notes whose encounters ends between May 1, 2021 and April 30, 2022, with additional H&P notes from January 2024 for the temporal test. A positive label is assigned if claim status is "final, adverse determination" (initial rejection by insurance and again rejected following appeal) or "final, favorable determination" (initial rejection by insurance and approved following appeal), an negative label otherwise. We split the dataset into five sets. The first three sets are train, validation, and test set with a 8:1:1 ratio from May 2021-Feb 2022. The 2022 temporal test include notes from March-Apr 2022. The 2024 temporal test set includes notes from January 2024. The positive class ratio range from 12.01% to 13.90%. The dataset contains 87,974 patients, 97,837 notes, and 89,147,715 words.

NYU+ Charlson Comorbidity Index (CCI). CCI is a standard score used to quantify a patient’s chronic illness based on their medical history [57]. It is useful for supporting accurate risk stratification for patients with limited historical data. However, this heuristic score cannot be computed when a patient’s medical history is unknown. This dataset addresses this problem by providing History & Physical (H&P) notes paired with their corresponding binned CCI scores. It allows models to be trained to impute the CCI score directly from unstructured clinical text. The target CCI is calculated using ICD code associated with the encounter in the EHR following the scoring function in [58]. Encounters with no associated ICD codes are excluded. The CCI is discretized into five classes: comorbidity index of 0 (<50th percentile), comorbidity index of 1 – 2 (50 – 75% percentile), comorbidity index of 3 – 4 (75-90th percentile), and comorbidity index of 5 – 7 (90-99th percentile), and comorbidity index >7 (>99th percentile). We split the dataset into five sets. The first three sets are train, validation, and test set with a 8:1:1 ratio from 2013 to May 2021. The 2021 temporal test includes notes from June to December 2021. The 2024 temporal test set includes notes from 2024. The majority class (score 0) ratio range from 68.40% to 69.47%. The minority class (score more than 7) ratio range from 0.059% to 0.17%. In total, the dataset has 306,741 patients, 443,915 notes, and 524,739,038 words.

5.2.3 External Validation Datasets

We also create external validation datasets from MIMIC III [59], which is sourced from Beth Israel hospital in Boston Massachusetts.

MIMIC III Readmission. The labeled dataset has 6% positive labels, with 52,725 examples and a 70% train, 15% validation and 15% test split. More dataset construction details are in [60]’s Appendix A.

MIMIC III Mortality. The labeled dataset has 10.55% positive labels, with 5658 examples and 80% train, 10% validation and 10% test split. The dataset is constructed from MIMIC-III clinical notes. As no explicit "Admission Note" label exists, we identify notes by filtering descriptions for "admission" while excluding "readmission" ($\approx 19K$ notes). To prevent patient bias, we select a single note per hospital stay using a prioritization heuristic. The heuristic first prioritizes Physician Note > General Note > Nurse Note, then prioritizes "resident/attending" > "resident" > "attending" > "fellow" > "h&p". Finally the heuristic prefers longer notes. This

filters down to ≈ 6 K notes. We further refine the dataset by removing notes written > 120 hours after admission or associated with a negative length of stay. The final dataset contains 5,658 unique admission notes.

MIMIC III LOS. The labeled dataset uses the same 5,658 admission notes as the mortality task, with a mean LOS of 7.96 days. Similarly, the split is 80% train, 10% validation and 10% test. The continuous LOS values are discretized into bins. We adapt the NYU+ LOS scheme (Methods 5.2.2) to handle MIMIC-III’s continuous range by treating the original integer bins as upper bounds. For example, the "3 days" bin is modified to capture the continuous range $2 < \text{LOS} \leq 3$.

5.2.4 Comprehension Datasets

We evaluate the performance of LANG1 checkpoints on comprehension datasets to analyze the emergence of its nonclinical abilities.

SciQ [61]. The dataset contains 13.7K multiple choice science exam questions with contexts. An example question is Mesophiles grow best in moderate temperature, typically between 25°C and 40°C (77°F and 104°F). Mesophiles are often found living in or on the bodies of humans or other animals. The optimal growth temperature of many pathogenic mesophiles is 37°C (98°F), the normal human body temperature. Mesophilic organisms have important uses in food preparation, including cheese, yogurt, beer and wine. \n Question: What type of organism is commonly used in preparation of foods such as cheese and yogurt? \n Answer:

PubmedQA [62]. The dataset contains 1k expert-annotated biomedical question answering dataset collected from PubMed abstracts. An example question is Abstract: Complex regional pain syndrome type I is treated symptomatically ... Early cast-related complaints predicted the development of complex regional pain syndrome (relative risk, 5.35; 95% confidence interval, 2.13 to 13.42)\n Question: Can vitamin C prevent complex regional pain syndrome in patients with wrist fractures?\n Answer:

5.3 Pretraining LANG1

We pretrain a family of Llama-style decoders (LANG1-100M, LANG1-1B, LANG1-7B) on a mixture of web texts and NYU Notes+ (section 5.2.1) using next token prediction (Figure 1 a,b). Detailed demographic statistics is in Appendix G. Unless otherwise noted, LANG1 models are trained with equal sampling from both clinical and general sources, which is supported by our pretraining ablations (Appendix D). For tokenization, we use the LLAMA-2-7B tokenizer (SentencePiece, 32k vocabulary). The 100M-parameter model follows the SMOL-LLAMA-101M architecture with a 1024 context length; the 1B model follows TINYLLAMA-1.1B with a 2048 context length; and the 7B model follows LLAMA-2-7B with a 4096 context length

We pretrain on 8 to 64 nVidia 80GB H100s with NVLink. We used the LitGPT [63] library and Fully Sharded Data Parallel [64].

We run a few trials of manual hyperparameter search based on speed, performance and training stability. For all models we use AdamW optimizer with linear warmup (2000 steps), beta1 of 0.9, beta2 of 0.95, eps of 1e-8 and cosine cycle decay down to a minimum learning rate of 4e-5. We use a seed of 3407, weight decay of 0.1, and gradient clipping of 1. In terms of FSDP sharding, we shard gradient and optimizer for models up to 1B, and full sharding for 7B model. For effective batch size, we use 4096 for 100M model for controlled comparison with NYUTron [42], and 1024 for 1B and 7B models.

We implement a monitoring pipeline that automatically triggered few-shot evaluations and generations at fixed intervals of pretraining steps. Slack Weights & Biases (W&B) alerts are configured to report loss spikes. Upon detection of anomalies in loss or output, we revert to the most recent stable checkpoint. Validation loss is computed periodically on the held-out 0.1% validation split and used to determine early stopping and model checkpoint selection.

5.3.1 Pretrained Models

We pretrain the following variants of LANG1 models (Table 1). Ablations (Appendix D) show that larger models trained on more clinical data has better performance, and mixing in web texts does not hurt much. When we refer to LANG1 without specifying data sources, we mean the one trained with NYUNotes+ and WebTexts. For instance, LANG1-1B refers to LANG1-1B-NYUNOTES+,WEBTEXTS.

Table 1: Pretrained Model Specifications

Model Name	Model Size	Pretrain Data
LANG1-100M-NYUNOTES	100m	NYUNotes
LANG1-100M-NYUNOTES+	100m	NYUNotes+
LANG1-100M-NYUNOTES+,WEBTEXTS	100m	NYUNotes+,WebTexts
LANG1-1B-NYUNOTES	1B	NYUNotes
LANG1-1B-NYUNOTES+	1B	NYUNotes+
LANG1-1B-NYUNOTES+,WEBTEXTS	1B	NYUNotes+,WebTexts
LANG1-7B-NYUNOTES+,WEBTEXTS	7B	NYUNotes+,WebTexts

5.4 Finetuning

We finetune LANG1 models (and their trajectories of checkpoints) and other pretrained models (Table 2) on ReMedE tasks using multiple choice format (Figure 1 panel c). The labeled clinical notes are converted to multiple choice format (Appendix H), and we train the model to predict the correct multiple choice option. For fair comparison with NYUTron, we right truncate all clinical notes to a maximum of 512 tokens. All finetuning jobs are done one node of 8 nVidia 80GB H100s with NVLink.

Before each full finetuning, we run 5 hyperparameter search trials up to 100 steps using Hydra and Optuna. We random search learning rate (based on [65]’s recommendation) in log scale over the closed interval 1e-6 to 1e-3. We used an AdamW optimizer with beta1 of 0.9, beta2 of 0.999, eps of 1e-5. We used a weight decay of 0.02 and no gradient clipping. We used a cosine annealing scheduler with no warmup and a max steps of 5120. The best trial is selected based on both maximum validation AUROC and minimum validation loss.

Table 2: Additional Model Specifications

Model Name	Model Size	Pretrain Data
LLAMA-3.2-1B	1B	Unnamed public mix (9T tokens)
LLAMA-2-7B	7B	Unnamed public mix (2T tokens)
DEEPSEEK-R1-DISTILL-LLAMA-70B	70B	Unnamed public mix (2T tokens) + reasoning data (800k samples)

For full finetuning, we used the best learning rate from hyperparameter search, and train for maximum 5120 steps with early stopping based on Micro-AUROC and a patience of 300 steps. The probabilities for calculating AUROC is obtained via normalizing the logits of the multiple choice options (e.g., “A”). We train all parameters except for DEEPSEEK-R1-DISTILL-LLAMA-70B, which has to be finetuned using low-rank adaptation finetuning [66] to meet the memory constraint (Appendix I).

For multitask finetuning, we mix examples from each task evenly within each training batch. We increase the total training steps scaled by the number of tasks.

5.5 Evaluation

Pretraining evaluation. We monitored the token-level cross entropy loss and perplexity for both training and validation.

Zero-shot and few-shot evaluation. ReMedE is built on LM Eval Harness [67]. We implemented the tasks as multiple choice questions and AUROC score as a metric. We implemented a child class of LocalCompletionsAPI to connect on-prem models. For models whose logits are not accessible (e.g., on-prem GPT-4o), we implemented custom sampling function to approximate (Appendix J) the probability via counting choices from 10 generations using a temperature of 1.

Finetuning evaluation. We collected the logits of the multiple choice options, normalized it as probabilities, and calculated AUROC using sklearn’s implementation (same as ReMedE’s AUROC backend for consistency). For multiclass classification, we used One-Versus-Rest (OVR) AUROC.

Uncertainty. To capture the uncertainty arose from the randomness of our test set, we calculated 95% confidence intervals ($CI = \pm 1.96 \times SD$) by resampling each test set 1000 times using the quantile bootstrap method from scipy.

Temporal Shift. To better mimic deployment conditions under temporal distribution shift, all AUROCs are reported on test data from 2024 – drawn from a period after the pretraining data – unless otherwise noted. See Appendix A for a visualization of the test timeline.

Generalist Models. We compared against generalist frontier models, including DEEPSEEK R1 (served via vLLM [68]), DEEPSEEK R1 DISTILLED LLAMA 70B (vLLM), LLAMA 3.3 70B CHAT (vLLM), and on-premises GPT-4o (Azure-hosted service). We also

evaluated MedQA leaderboard models using in-memory inference, such as LLAMA 3.2 1B, LLAMA 2 7B, and MEDMOBILE.

Acknowledgements. E.K.O. is supported by the National Cancer Institute’s Early Surgeon Scientist Program (3P30CA016087-41S1) and the W.M. Keck Foundation. L.Y.J. is supported by Apple AIML PhD fellowship. L.Y.J. and A.C. are supported by NSF Award 1922658. K.C., E.K.O., L.Y.J. and A.C. are supported by Institute for Information & communications Technology Promotion (IITP) grant funded by the Korea government (MSIT) (No. RS-2019-II190075 Artificial Intelligence Graduate School Program (KAIST); No. RS-2024-00509279, Global AI Frontier Lab). We would like to acknowledge J. Golfinos, whose vision and support made this project possible. We would like to acknowledge Michael Costantino, Ph.D., Ali Siavosh-Haghighi, Ph.D., Kevin Yie, M.S., Neelima Sharma, Tedum Sampson from the NYU Langone High Performance Computing (HPC) team. Without their tireless assistance in building and maintaining our GPU cluster none of this research would have been possible. We would also like to thank Dr. Dafna Bar-Sagi, Ph.D., and Nader Mherabi whose support for this research has made everything possible. Thanks to He He, Ph.D., Eunsol Choi, Ph.D., Carlos Fernandez-Granda, Ph.D., Julia Kempe, Ph.D., Vasant Dhar, Ph.D., Keunwoo Choi, Ph.D., Jesse Swanson, Gavin Zihao Yang, William Merrill, Ph.D., Nicholas Lourie, Sophie Hao, Ph.D., Vishakh Padmakumar, Ph.D., Michael Hu, Robert J Steele, Yueying Li, Yunzhen Feng, Guillermo Sapiro, Ph.D., Oussama Elaqqhar, Kai Xu, Varun Yerram, Itay Itzhak, Jeff Hammerbacher for their valuable discussions.

Declarations

Ethical approval. This study was approved by the Institutional Review Board (IRB) at NYU Langone Health (study protocol s21-01189). The methods were carried out in accordance with the IRB’s relevant guidelines and regulations.

Data Availability. The clinical data used for the pretraining, finetuning, validation, and test sets were collected from the NYU Langone Health System EHR maintained by the NYULH Datacore team. Text data was stripped of rich text features and directly included in the dataset "as-is", and was augmented with structured features where noted. It consists of the production medical records of NYU Langone and cannot be made publicly available. For the external validation task, the datasets were obtained from MIMIC III, and are publicly available from their [website](#).

Code Availability. This work uses several open-source libraries including [PyTorch](#), [LitGPT](#), [Transformers library](#), [LM Eval Harness](#), and [hydra](#). Our experimental framework involves the utilization of these libraries and in some cases modification of them. We will release code to replicate the pretraining, finetuning, and testing of the models described in this paper at the time of publication. We include detailed methods and implementation steps in the Methods and Supplementary Information to allow for independent replication.

Author’s contribution. E.K.O. and K.C. supervised the project. L.Y.J. and X.H. collected pretrain, finetune and evaluation data (except NYU+ Insurance Denial and MIMIC-III). L.Y.J. and A.C. engineered the software for pretrain and finetune. A.C., X.H., L.Y.J. and J.Z. engineered the software for evaluation. L.Y.J., A.C., X.H., R.D., X.C.L. ran experiments. L.Y.J., A.C., X.C.L., X.H., J.Z., R.D. and K.C. debugged and tested the software. A.C., L.Y.J., K.C.,

X.C.L., R.S., A.A., K.L.S, Y.C., and Q.P. created figures. A.C., L.Y.J., K.C., E.K.O. conceptualized the training dynamics and transfer experiments. J.S. hosted deepseek and llama 70b inference server. K.E. collected NYU+ Insurance Denial data. Q.P. and L.Y.J. check pretrain data quality and cleaned the data. F.W. processed the MIMIC-III Mortality and LOS data. A.C., L.Y.J, R.D., Y.C., D.A. and J.V.L. reviewed related literature. K.C., E.K.O., Y.A. provided guidance and feedback throughout the project. L.Y.J., A.C., X.H., R.D., A.A., D.A., J.V.L., Q.P., Y.C., R.J.S., K.C., E.K.O. wrote the initial draft. All authors edited and revised the manuscript.

Conflict of interest E.K.O. reports consulting with March AI, Sofinnova Inc., Google Inc., income from Merck & Co., and Mirati Therapeutics, and equity in Artisight Inc. K.C. is employed by Prescient Design, a subsidiary of Roche. A.C. is employed by Google Deepmind. Q.P. is employed by Faction Imaging Inc. J.S. is employed by March AI. There are no other potential conflicts of interest. The work presented herein was performed exclusively within the NYU Langone Health System.

Appendix A Data Timeline

LANG1’s pretrain data covers a wider time window than NYUTRON[42]. NYUTRON’s pretrain data is from 2013 to May 2021. LANG1’s pretrain data covers a longer span, from 2003 to 2023. Overall, LANG1’s pretrain corpus is more than 10 times the size of NYUTRON.

LANG1 uses the same finetuning dataset for training as NYUTRON, but adds additional temporal test. For both NYUTRON and LANG1 finetuning, we have a temporal test set to better mimic the deployment scenario, where the test set comes from the future of training set. For NYUTRON, the temporal test set is between June to December of 2021. For LANG1, the temporal test set is in 2024. All performance we report in main is performance on 2024 temporal test set, unless otherwise specified.

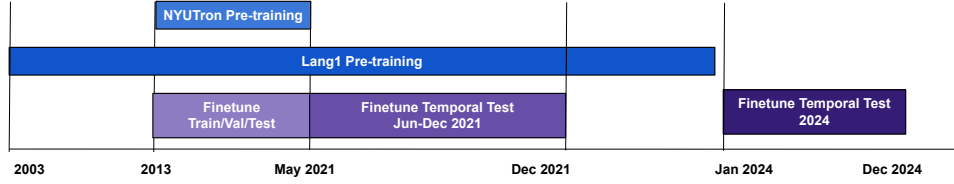


Fig. A1: Illustration of timeline of pretrain and finetune dataset for both LANG1 and NYUTRON. LANG1 covers a wider time window for pretraining and added additional temporal test set in 2024 to capture temporal distribution shift.

Temporal test is important and difficult. Figure A2 shows that both LANG1-1B (purple) and NYUTRON (pink) perform worse as test data are sampled from a further future, illustrating the importance of evaluating with temporal test to capture deployment scenario.

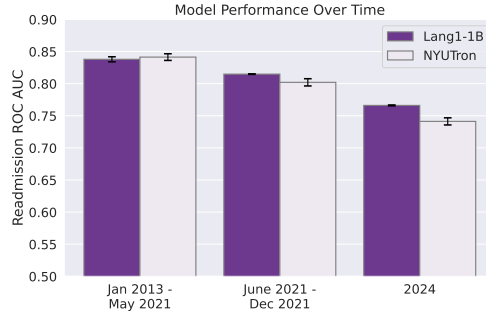


Fig. A2: Models perform worse on temporal test and exhibit different level of degradation in face of temporal shift.

Appendix B ReMedE Dataset Statistics

The following tables show the note counts (Table B1) and patient counts (Table B2) across each split, and class ratio for each task: readmission (Table B3), mortality (Table B4), length of stay (Table B5, comorbidity imputation (Table B6, and Insurance denial (Table B7).

Table B1: Distinct note counts for each ReMedE task across five splits.

Task	Train	Val	Test	Temporal Test 2021	Temporal Test 2024	Total
Readmission	362,259	45,282	45,283	53,916	97,586	604,326
CCI	256,676	32,085	32,085	42,137	80,932	443,915
Length of Stay	334,515	41,814	41,815	51,018	97,586	566,748
Insurance Denial	41,842	2,325	2,325	9,299	42,046	97,837
Mortality	334,515	41,814	41,815	51,018	97,586	566,748

Table B2: Distinct patient counts for each ReMedE task across five splits.

Task	Train	Val	Test	Temporal Test 2021	Temporal Test 2024	Total
Readmission	269,140	42,692	42,603	46,003	78,453	421,429
CCI	188,298	30,085	30,098	251,804	64,873	306,741
Length of Stay	248,486	39,304	39,331	43,358	78,453	395,991
Insurance Denial	39,422	2,319	2,313	9,037	37,821	87,974
Mortality	248,674	39,317	39,357	43,358	78,453	395,991

Table B3: Readmission label ratios by split (Total notes = 604,326; total words = 607,877,177).

Split	Not Readmitted	Readmitted within 30 days
Train	0.891495	0.108505
Val	0.891944	0.108056
Test	0.893293	0.106707
Temporal Test (2021)	0.888456	0.111544
New Temporal (2024)	0.887115	0.112885

Table B4: In-hospital mortality label ratios by split (Total notes = 566,748; total words = 608,603,182).

Split	Survived	Died
Train	0.981161	0.018839
Val	0.981250	0.018750
Test	0.980916	0.019084
Temporal Test (2021)	0.980693	0.019307
New Temporal (2024)	0.982241	0.017759

Table B5: Length of stay label ratios by split (Total notes = 566,748; total words = 608,603,182).

Split	0–2 days	3 days	4–5 days	>5 days
Train	0.417306	0.176575	0.165888	0.240231
Val	0.415722	0.180562	0.164490	0.239226
Test	0.414851	0.176611	0.167452	0.241086
Temporal Test (2021)	0.418597	0.153201	0.164805	0.263397
New Temporal (2024)	0.456418	0.144529	0.159757	0.239297

Table B6: Charlson Comorbidity Index (CCI) label ratios by split (Total notes = 443,915; total words = 524,739,038).

Split	Score 0	Score 1–2	Score 3–4	Score 5–7	Score >7
Train	0.694681	0.224840	0.054251	0.025324	0.000904
Val	0.689169	0.229547	0.055166	0.025526	0.000592
Test	0.693502	0.226866	0.053327	0.025370	0.000935
Temporal Test (2021)	0.685811	0.228327	0.059164	0.025821	0.000878
New Temporal (2024)	0.684043	0.217355	0.067513	0.029370	0.001717

Table B7: Insurance denial label ratios by split (Total notes = 97,837; total words = 89,147,715).

Split	Approved (0)	Denied (1)
Train	0.877850	0.122150
Val	0.867097	0.132903
Test	0.873978	0.126022
Temporal Test (2021)	0.879880	0.120120
New Temporal (2024)	0.861033	0.138967

Appendix C Transfer Pattern of LLAMA 3.2 1B v.s. LANG1 1B

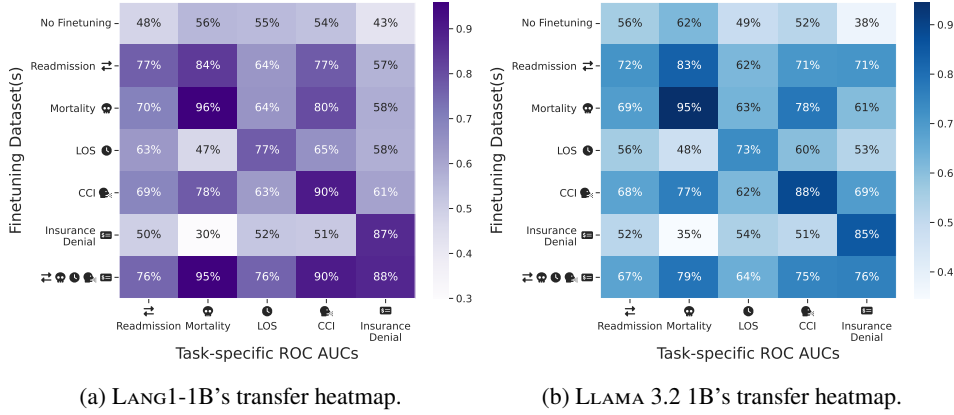


Fig. C3: LANG1-1B and LLAMA 3.2 1B transfers differently. The heatmap shows the two different base models' performance when finetuned on a subset of REMeDE tasks (y axis) and evaluated on all five REMeDE tasks (x axis). Overall, LANG1-1B has higher per-task and joint performance, and shows a different transfer pattern than LLAMA 3.2 1B.

LLAMA-3.2-1B has overall worse performance than LANG1-1B. Compared to Figure C3a, Figure C3b has worse single-task (diagonal) and joint-task (last row) performance, suggesting that the specific transfer pattern could be highly model specific.

Both models exhibit some similar patterns. (i) finetuning on readmission (second row) boosts performance on the other four tasks, and (ii) transfer could be asymmetric, *i.e.*, mortality helps LOS (third row, third column) but LOS does not help mortality (fourth row, second column), which can be explained using domain knowledge (Appendix F).

LANG1-1B and LLAMA-3.2-1B also have different patterns. (i) joint finetuning (last row) helps LANG1-1B but hurts LLAMA-3.2-1B, and (ii) finetuning on insurance denial (fourth row) lowers LANG1-1B's LOS performance (third column) while improving it for LLAMA-3.2-1B. These results suggest that **instruction finetuning enables cross-task transfer, though the specific transfer patterns depend on model pretraining.**

Appendix D Pretraining ablations

Figure D4 presents pretraining ablation results evaluated on 2024 readmission temporal test set. For pretraining, we control for the model architecture (encoder v.s. decoder), model size (100M v.s. 1B), and pretrain data (NYUNotes, NYUNotes+, NYUNotes+ and web texts). For finetuning, we evaluated on three clinical predictive tasks: readmission prediction, insurance denial prediction, and LOS prediction. These three tasks were chosen for their distinct transfer pattern in Figure C3a.

Training larger models on more recent clinical data improves temporal robustness. Figure D4a shows the ablation results for readmission prediction. On the left, holding the model architecture fixed as a decoder, we vary pretraining data. Compared to models trained only on EHR from 2013 to 2021 ($\mathcal{D}_{\text{NYUNotes}}$), adding more recent clinical data ($\mathcal{D}_{\text{NYUNotes+}}$, spanning 2003 to 2023) and further mixing in general-domain SlimPajama ($\mathcal{D}_{\text{NYUNotes+}, \text{WebText}}$) improves performance for the 1B model but not the 100M model, suggesting that larger models are better able to leverage additional clinical data. **This also justifies our choice of equally mixing EHR texts and web texts**, since adding web texts does not significantly hurt the downstream clinical task performance while instilling general-purpose knowledge. On the right, holding the pretraining data fixed to $\mathcal{D}_{\text{NYUNotes}}$, varying model architecture shows that scaling to 1B helps. We observe similar patterns for insurance denial prediction (Figure D4b) and LOS prediction (Figure D4c).

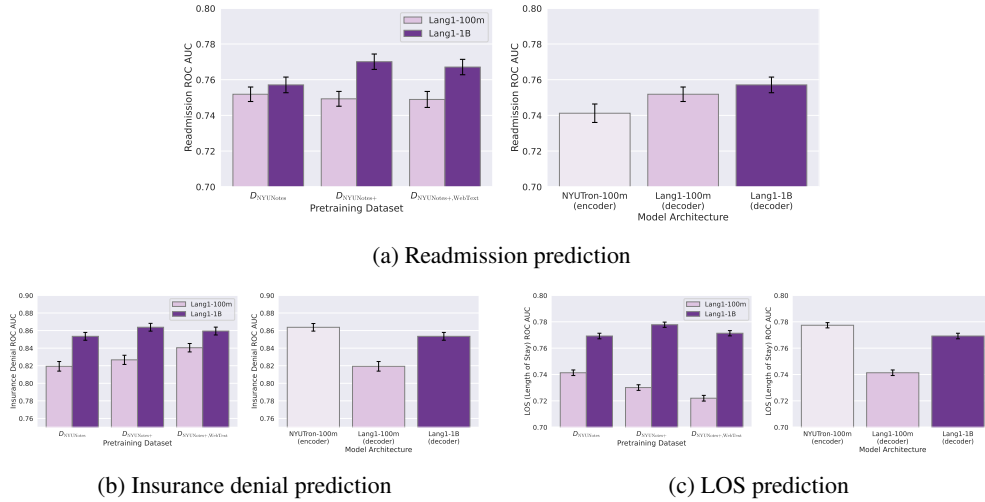
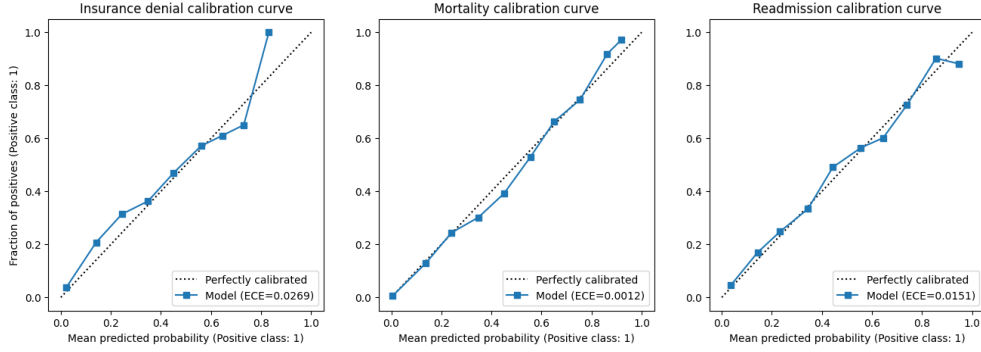


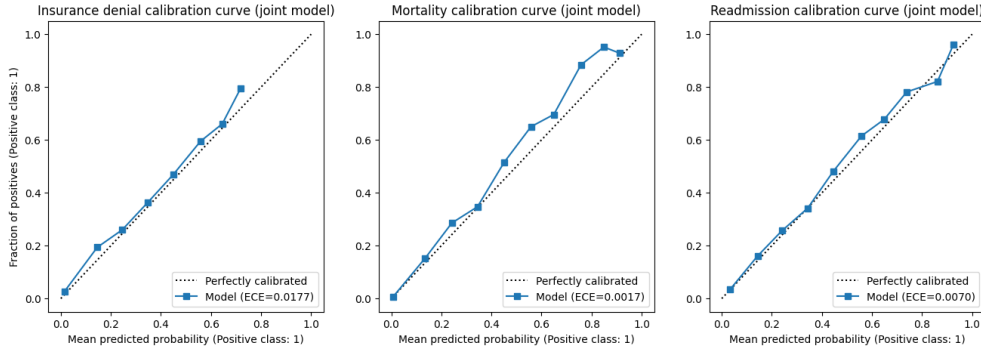
Fig. D4: Pretraining ablations. Error bars indicating the 95% confidence interval. (a) Larger models trained on more clinical data has better performance, and mixing web texts does not hurt much. (b,c) Similar patterns are observed for insurance denial and LOS prediction.

Appendix E Calibration Plot

Calibration curves are calculated using sklearn package [69] with $n=15$ bins. Expected calibration error (ECE) is calculated with $n=15$ bins using the torchmetrics library [70].



(a) Calibration plots for single-task models



(b) Calibration plots for joint model

Fig. E5: Calibration plot shows that both single task finetuned Lang1 and joint finetuned Lang1 are well calibrated.

Appendix F Asymmetry of Transfer between Mortality and LOS

Our medical collaborators provided an explanation for the asymmetric transfer between Mortality and LOS (Length of Stay) we observe for both LANG1-1B and LLAMA-3.2-1B. If a patient died, they either stayed for a short time (very sick and died immediately) or a long time (doctors failed to save them after a long time). On the other hand, if a patient stayed for a long time, they either survived, or they died after doctors' attempts. This asymmetry in conditional probability could help explain why the mortality transfer to LOS, but not vice versa.

This explanation is corroborated by analysis of the conditional probabilities. Figure F6a shows that patient who died are ore likely to stay for 0-2 days or >5 days compared to patient who survived. F6b shows that patient who stay for a long time have similar mortality risks as patient who stayed for 0 days.

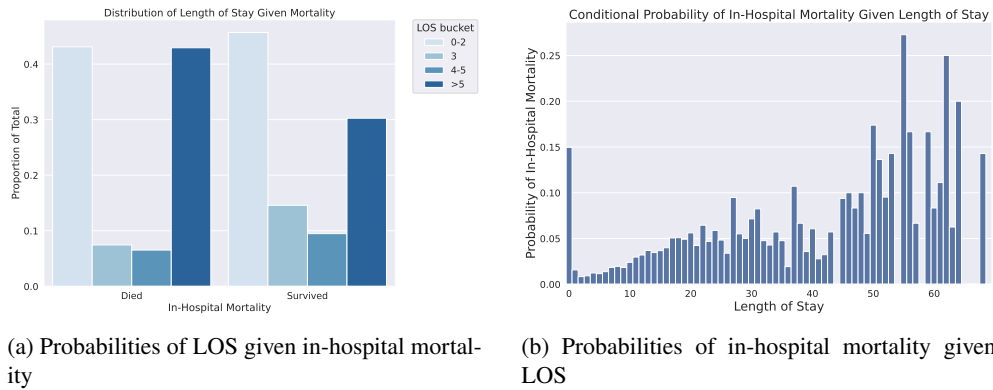
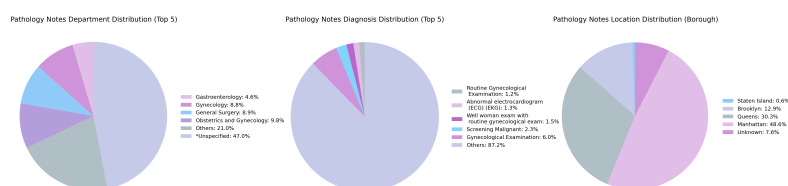


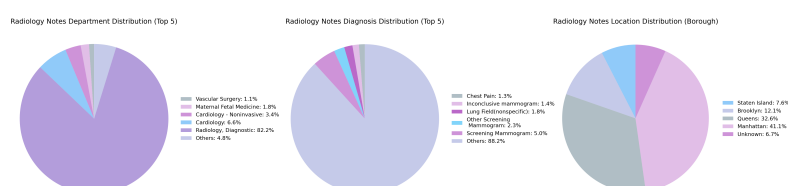
Fig. F6: Analysis of conditional probabilities of LOS (length of stay) and in-hospital mortality could help explain the assymetry of transfer.

Appendix G Detailed statistics of NYU Notes+

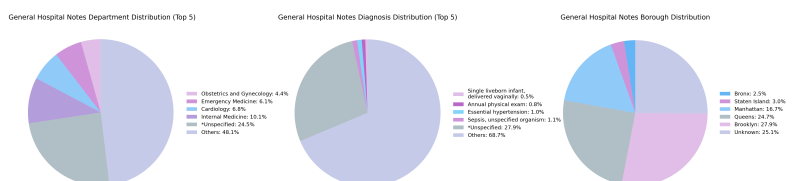
We analyzed the top five clinical departments, diagnosis and borough for pathology notes, radiology notes and hospital notes. See Figure G7.



(a) **Pathology Notes.** Among pathology notes with specified departments, OB/GYN (obstetrics and gynecology) has the highest percentage (9.8%). The most common specified diagnosis is gynecological exams (6%). Nearly half (48.6%) of the pathology notes are from Manhattan borough.



(b) **Radiology Notes.** Most of the radiology notes are from diagnostic radiology (82.2%), with screening mammograms being the most common specified diagnosis (5%). The two most common boroughs are Manhattan (41%) and Queens (32.6%).



(c) **Hospital Notes.** Among notes with specified departments, most are from internal medicine department (10.1%), with common diagnoses including sepsis (1.1%) and hypertension (1%). The majority of notes are from Brooklyn (27.9%) and Queens (24.7%) borough.

Fig. G7: Top five department, diagnosis and borough of NYU Notes+

Appendix H Prompts of ReMedE Tasks

We constructed prompts to create task-specific questions and answer options from the labeled finetuning notes:

- **Readmission** Question: Given the above discharge note of the patient, will the patient be readmitted to the hospital within 30 days of discharge? \n A. no \n B. yes \n Answer:
- **In-Hospital Mortality** Question: Given the above admission note of the patient, will the patient die during the hospital admission? \n A. no \n B. yes \n Answer:
- **Charlson Comorbidity Index** Question: Given the above admission note of the patient, what's the Charlson Comorbidity Index of the patient? \n A. score 0 \n B. score 1 to 2 \n C. score 3 to 4 \n D. score 4 to 7 \n E. score more than 7 \n Answer:
- **Insurance Denial** Question: Given the above discharge note of the patient, will the insurance claim of the patient be denied? \n A. no \n B. yes \n Answer:
- **Length of Stay** Question: Given the above admission note of the patient, how long will the patient stay at the hospital? \n A. 0 to 2 days \n B. 3 days \n C. 4 to 5 days \n D. more than 5 days \n Answer:

Appendix I LoRa Finetuning for Llama-3-70b

To efficiently train Deepseek-R1-distill-Llama-3-70b on 1 node of 8 H100s, we used Low-Rank Approximation (LoRA) finetuning. LoRA reduces trainable parameters by inserting trainable rank decomposition matrices into transformer layers while freezing the pretrained weights.

In our configuration, we enabled LoRA adapters on the query and value projections of the attention mechanism, with rank $r = 8$. We set the LoRA scaling factor $\alpha = 16$ and applied a dropout rate of 0.05. Other components, such as the key, MLP, and projection layers, were left frozen.

Appendix J Token probability approximation for models without logprobs

Unlike other generalist models, GPT-4o does not provide logprobs to prevent privacy attack. However, we need probabilities to reliably evaluate classification tasks. To approximate its probabilities, we sample 10 generations from GPT using temperature of 1 (to not reweight the LM's distribution), and count the number of occurrences for each multiple choice options. Then we normalize the counts to be probabilities for each option. For cost reasons, we limit the number of examples to be 1000, except for CCI (comorbidity imputation).

CCI is a special case because its label distribution is very skewed, so we greedily evaluate on 10,000 samples instead. We binarize the greedy choice to be probability of 0 or 1.

Appendix K Stratified Evaluation

We performed stratified evaluation on readmission to evaluation the performance variation across different groups (age, first race, borough, ethnicity, sex, and whether the patients are children). Some groups are omitted because they have only one class due to small sample size. If we look at means only, for age (Figure K8a), patients between 10 to 15 has the best performance, and patients between 85 and 90 has the worst performance. For race, Middle Eastern or North African has the best performance, and Asian Indian has the worst performance. For borough, Brooklyn has the best performance and Queens has the worst performance. For ethnicity, Spanish Hispanic Origin is worst. For sex, female has better performance than male. Children's performance is better than adult and they also have a lower readmission rate.

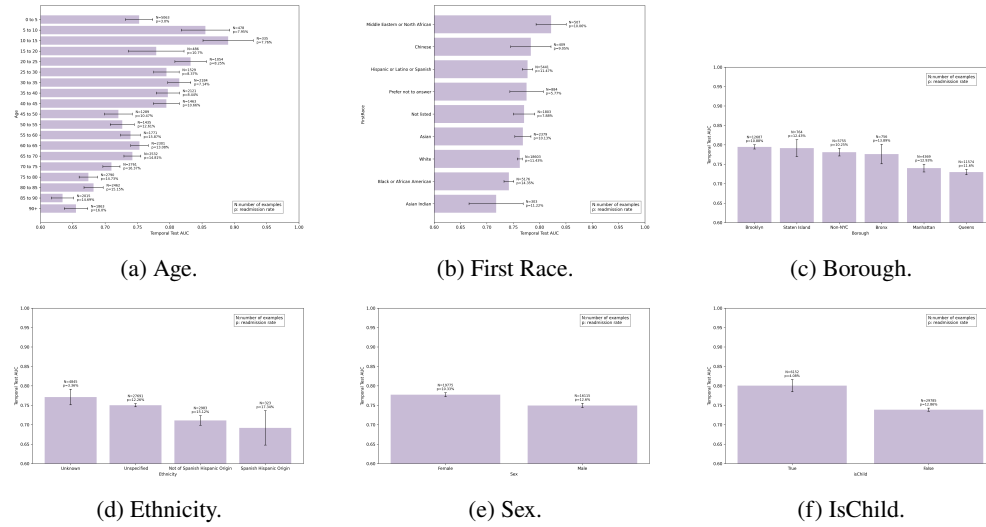


Fig. K8: Stratified performance analysis on 2024 readmission temporal test set.

Appendix L Control for Patient Overlap

We chose to construct the 2024 temporal test set without explicit patient split, because patients do come back to the health system at deployment. While this choice is supported by our clinical collaborators, we did an ablation where the test set excludes seen patients. Keeping the finetuned model fixed and varying the test data to either include or exclude seen patients, we see **similar test performance on readmission prediction, in hospital mortality prediction, and length of stay prediction**. Figure L9a shows that on readmission prediction, removing patients seen from pretraining and finetuning slightly increases the performance from 76.5% (purple bar) to 77% (grey bar). Similarly, the performance increase is 0.39% for mortality (Figure L9b) and 0.70% for LOS (Figure L9c). While surprising, the slight increase could be attributed to **repeated patients being both older and more likely to be a minority**. For readmission prediction, repeated patients are on average 13 years older than non repeated patients (55 v.s. 42 years). For mortality and LOS prediction, the age gap widens to 16 years (57 v.s. 41 years). In addition, repeated patients include a larger proportion of non-white patients (41% v.s. 36%). These findings show that our choice of temporal split does not over-estimate model performance.

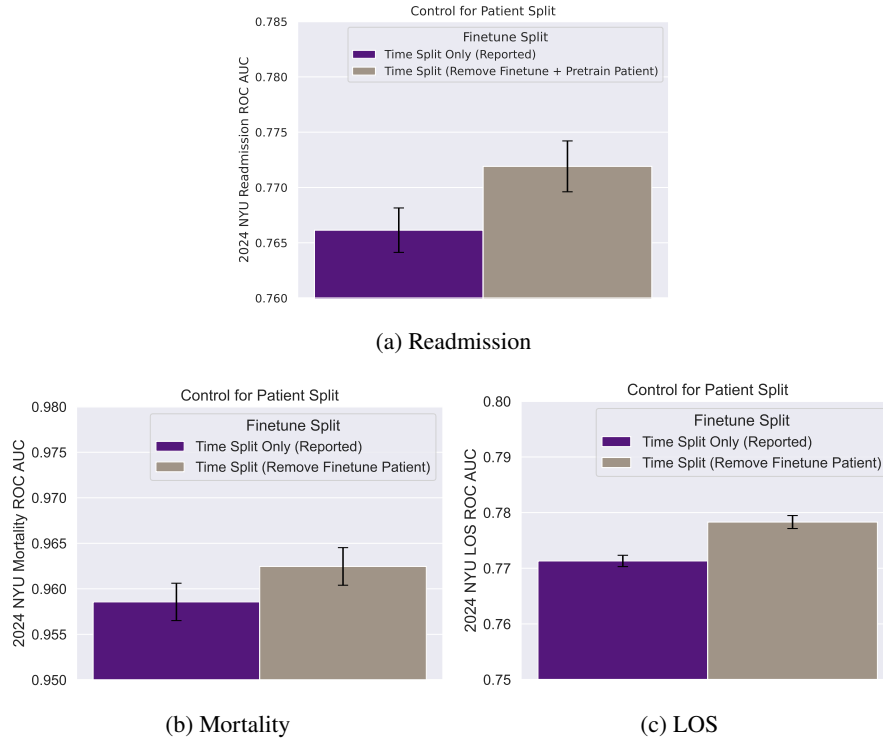
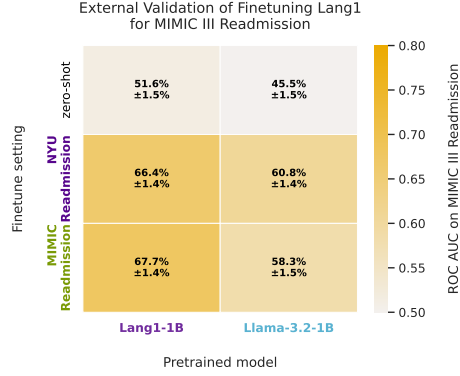


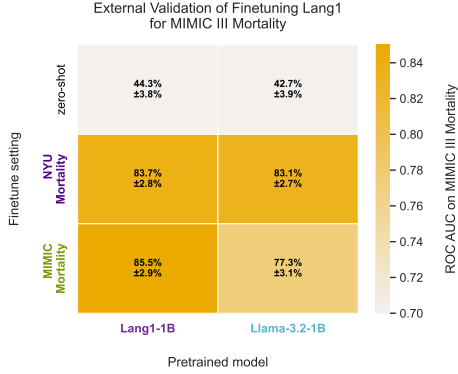
Fig. L9: LANG1-1B’s finetuned performance on readmission, mortality and LOS 2024 test set, with or without patient split.

Appendix M External Validation: MIMIC v.s. NYU

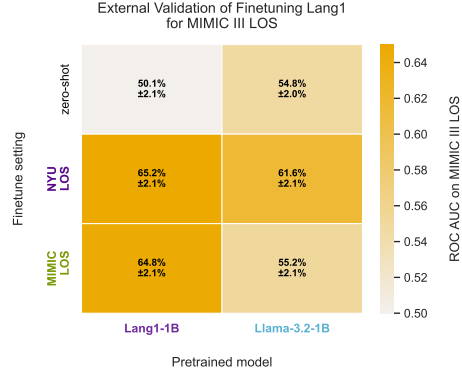
To check how well LANG1 generalizes to a different health system, we compare finetuning LANG1-1B and LLAMA-3.2-1B on three classification tasks (readmission, mortality, length of stay) from MIMIC III and NYU, and test on MIMIC III. MIMIC III [59] is a database of electronic health records from ICU patients at the Beth Israel Hospital in Boston Massachusetts, whereas NYU+ datasets are collected from New York City. See 5.2.3 for details about MIMIC datasets and Methods 5.2.2 for NYU datasets. Figure M10 shows the heatmaps of MIMIC III test performance for finetuning LANG1-1B and LLAMA-3.2-1B on NYU or MIMIC-III, with zero-shot performance as baselines. The x axis is the pretrained model (purple LANG1-1B and blue LLAMA-3.2-1B). The y axis is the finetuning setting (purple indicates finetuning on NYU data; green indicates finetuning on MIMIC III data; and black indicates zero-shot). Darker yellow cells indicate higher test AUROC on MIMIC III.



(a) Readmission prediction.



(b) In-hospital Mortality.



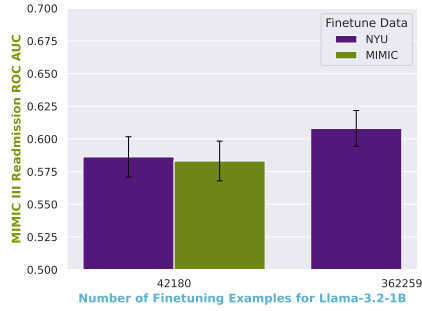
(c) LOS.

Fig. M10: Finetuning on NYU transfers well to MIMIC III for readmission, mortality and LOS prediction. Overall performance is better on finetuning LANG1-1B compared to LLAMA 3.2 1B.

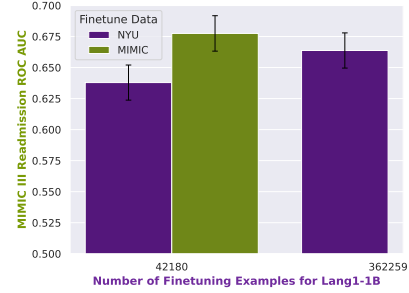
Finetuning LANG1-1B on NYU Readmission transfers to MIMIC-III Readmission. The difference in mean AUROC between finetuning LANG1-1B on MIMIC III versus NYU range from 0.4%-1.8% AUROC and are roughly within standard error, showing that finetuning LANG1-1B at NYU transfers well to MIMIC.

LANG1-1B achieves better performance than LLAMA-3.2-1B. The difference in mean AUROC between finetuning LANG1-1B and LLAMA-3.2-1B on the same data range from 0.6%-9.6%, showing that the clinically pretrained LANG1-1B is a preferred base model for these three clinical predictive tasks.

Finetuning on NYU, which is slightly out-of-distribution, could outperform finetuning on MIMIC for LLAMA-3.2-1B. It is reasonable to expect that in-distribution finetuning leads to best performance. While this is true for LANG1-1B (purple), it is not the case for LLAMA-3.2-1B (blue), which sees 2.5% - 6.4% AUROC increase from finetuning on NYU compared to finetuning on MIMIC. We hypothesize that this is because NYU data has more labeled pairs, and that **non clinically pretrained models would benefit more from a larger number of slightly out-of-distribution pairs**. We test this hypothesis on readmission prediction, where NYU Readmission (36,2259 examples) is 8.6 times the size of MIMIC III readmission (42,180 examples). Figure M11a shows that test performance on MIMIC III is similar when LLAMA-3.2-1B is finetuned on NYU Readmission that is **downsampled** to be the same size as MIMIC III (42,180). This pattern does not hold when the pretrained model is changed to the clinically pretrained LANG1-1B: Figure M11b shows that finetuning on MIMIC-III readmission yields the best test performance on MIMIC-III, despite fewer number of finetuning examples.



(a) Finetuning on full NYU dataset leads to best performance on LLAMA-3.2-1B, even though NYU Readmission is not directly in-distribution for MIMIC Readmission.



(b) Finetuning on MIMIC III leads to best performance for LANG1-1B, even though NYU Readmission has more labeled pairs.

Fig. M11: Clinical models such as LANG1-1B benefit more from in-domain data, whereas generalist models such as Llama-3.2-1B could benefit from more slightly OOD examples compared to fewer in-distribution examples.

References

- [1] Sinsky, C., Colligan, L., Li, L., Prgomet, M., Reynolds, S., Goeders, L., Westbrook, J., Tutty, M., Blike, G.: Allocation of physician time in ambulatory practice: A time and motion study in 4 specialties. *Ann. Intern. Med.* **165**(11), 753–760 (2016)
- [2] Hingle, S.: Electronic health records: An unfulfilled promise and a call to action. *Ann. Intern. Med.* **165**(11), 818–819 (2016)
- [3] Murphy, D.R., Meyer, A.N.D., Russo, E., Sittig, D.F., Wei, L., Singh, H.: The burden of inbox notifications in commercial electronic health records. *JAMA Intern. Med.* **176**(4), 559–560 (2016)
- [4] Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., Cao, Y.: ReAct: Synergizing reasoning and acting in language models. *arXiv [cs.CL]* (2022) [cs.CL]
- [5] Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., Casas, D., Hendricks, L.A., Welbl, J., Clark, A., Hennigan, T., Noland, E., Millican, K., Driessche, G., Damoc, B., Guy, A., Osindero, S., Simonyan, K., Elsen, E., Vinyals, O., Rae, J.W., Sifre, L.: An empirical analysis of compute-optimal large language model training. In: Oh, A.H., Agarwal, A., Belgrave, D., Cho, K. (eds.) *Advances in Neural Information Processing Systems* (2022). <https://openreview.net/forum?id=iBBcRUIOAPR>
- [6] DeepSeek-AI, Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., Zhang, X., Yu, X., Wu, Y., Wu, Z.F., Gou, Z., Shao, Z., Li, Z., Gao, Z., Liu, A., Xue, B., Wang, B., Wu, B., Feng, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., Dai, D., Chen, D., Ji, D., Li, E., Lin, F., Dai, F., Luo, F., Hao, G., Chen, G., Li, G., Zhang, H., Bao, H., Xu, H., Wang, H., Ding, H., Xin, H., Gao, H., Qu, H., Li, H., Guo, J., Li, J., Wang, J., Chen, J., Yuan, J., Qiu, J., Li, J., Cai, J.L., Ni, J., Liang, J., Chen, J., Dong, K., Hu, K., Gao, K., Guan, K., Huang, K., Yu, K., Wang, L., Zhang, L., Zhao, L., Wang, L., Zhang, L., Xu, L., Xia, L., Zhang, M., Zhang, M., Tang, M., Li, M., Wang, M., Li, M., Tian, N., Huang, P., Zhang, P., Wang, Q., Chen, Q., Du, Q., Ge, R., Zhang, R., Pan, R., Wang, R., Chen, R.J., Jin, R.L., Chen, R., Lu, S., Zhou, S., Chen, S., Ye, S., Wang, S., Yu, S., Zhou, S., Pan, S., Li, S.S., Zhou, S., Wu, S., Ye, S., Yun, T., Pei, T., Sun, T., Wang, T., Zeng, W., Zhao, W., Liu, W., Liang, W., Gao, W., Yu, W., Zhang, W., Xiao, W.L., An, W., Liu, X., Wang, X., Chen, X., Nie, X., Cheng, X., Liu, X., Xie, X., Liu, X., Yang, X., Li, X., Su, X., Lin, X., Li, X.Q., Jin, X., Shen, X., Chen, X., Sun, X., Wang, X., Song, X., Zhou, X., Wang, X., Shan, X., Li, Y.K., Wang, Y.Q., Wei, Y.X., Zhang, Y., Xu, Y., Li, Y., Zhao, Y., Sun, Y., Wang, Y., Yu, Y., Zhang, Y., Shi, Y., Xiong, Y., He, Y., Piao, Y., Wang, Y., Tan, Y., Ma, Y., Liu, Y., Guo, Y., Ou, Y., Wang, Y., Gong, Y., Zou, Y., He, Y., Xiong, Y., Luo, Y., You, Y., Liu, Y., Zhou, Y., Zhu, Y.X., Xu, Y., Huang, Y., Li, Y., Zheng, Y., Zhu, Y., Ma, Y., Tang, Y., Zha, Y., Yan, Y., Ren, Z.Z., Ren, Z., Sha, Z., Fu, Z., Xu, Z., Xie, Z., Zhang, Z., Hao, Z., Ma, Z., Yan, Z., Wu, Z., Gu, Z., Zhu, Z., Liu, Z., Li, Z., Xie, Z., Song, Z., Pan, Z., Huang, Z., Xu, Z., Zhang, Z., Zhang, Z.: DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning (2025). <http://arxiv.org/abs/2501.12948>

- [7] Trinh, T.H., Wu, Y., Le, Q.V., He, H., Luong, T.: Solving olympiad geometry without human demonstrations. *Nature* **625**(7995), 476–482 (2024)
- [8] Gottweis, J., Weng, W.-H., Daryin, A., Tu, T., Palepu, A., Sirkovic, P., Myaskovsky, A., Weissenberger, F., Rong, K., Tanno, R., Saab, K., Popovici, D., Blum, J., Zhang, F., Chou, K., Hassidim, A., Gokturk, B., Vahdat, A., Kohli, P., Matias, Y., Carroll, A., Kulkarni, K., Tomasev, N., Guan, Y., Dhillon, V., Vaishnav, E.D., Lee, B., Costa, T.R.D., Penad s, J.R., Peltz, G., Xu, Y., Pawlosky, A., Karthikesalingam, A., Natarajan, V.: Towards an AI co-scientist (2025). <http://arxiv.org/abs/2502.18864>
- [9] Wu, S., Irsoy, O., Lu, S., Dabrowski, V., Dredze, M., Gehrmann, S., Kambadur, P., Rosenberg, D., Mann, G.: BloombergGPT: A Large Language Model for Finance (2023). <http://arxiv.org/abs/2303.17564>
- [10] Chen, J., Li, D., Chen, Q., Zhou, W., Liu, X.: Diaformer: Automatic diagnosis via symptoms sequence generation. *Proc. Conf. AAAI Artif. Intell.* **36**(4), 4432–4440 (2022)
- [11] Chen, W., Zhong, C., Peng, J., Wei, Z.: DxFormer: a decoupled automatic diagnostic system based on decoder-encoder transformer with dense symptom representations. *Bioinformatics* **39**(1), 744 (2023)
- [12] Panagoulas, D.P., Palamidis, F.A., Virvou, M., Tsihrintzis, G.A.: Evaluating the potential of LLMs and ChatGPT on medical diagnosis and treatment. In: 2023 14th International Conference on Information, Intelligence, Systems & Applications (IISA), pp. 1–9. IEEE, ??? (2023)
- [13] Amin, K., Khosla, P., Doshi, R., Chheang, S., Forman, H.P.: Artificial intelligence to improve patient understanding of radiology reports. *Yale J. Biol. Med.* **96**(3), 407–417 (2023)
- [14] Ellershaw, S., Tomlinson, C., Burton, O.E., Frost, T., Hanrahan, J.G., Khan, D.Z., Horsfall, H.L., Little, M., Malgapo, E., Starup-Hansen, J., Ross, J., Noor, K., Vella-Baldacchino, M., Shah, A.D., Dobson, R.: Automated generation of hospital discharge summaries using clinical guidelines and large language models. In: AAAI 2024 SSS on Clinical FMs, pp. 1–8. Stanford University, Stanford, California, USA (2024)
- [15] Zhou, S., Xu, Z., Zhang, M., Xu, C., Guo, Y., Zhan, Z., Ding, S., Wang, J., Xu, K., Fang, Y., Xia, L., Yeung, J., Zha, D., Melton, G.B., Lin, M., Zhang, R.: Large language models for disease diagnosis: A scoping review. *arXiv [cs.CL]* (2024) [cs.CL]
- [16] Zhang, L., Liu, M., Wang, L., Zhang, Y., Xu, X., Pan, Z., Feng, Y., Zhao, J., Zhang, L., Yao, G., Chen, X., Xie, X.: Constructing a large language model to generate impressions from findings in radiology reports. *Radiology* **312**(3), 240885 (2024)
- [17] He, Z., Bhasuran, B., Jin, Q., Tian, S., Hanna, K., Shavor, C., Arguello, L.G., Murray, P., Lu, Z.: Quality of answers of generative large language models versus peer users for interpreting laboratory test results for lay patients: Evaluation study. *J. Med. Internet*

- [18] Kweon, S., Kim, J., Kwak, H., Cha, D., Yoon, H., Kim, K., Yang, J., Won, S., Choi, E.: EHRNoteQA: An LLM benchmark for real-world clinical practice using discharge summaries. *Neural Inf Process Syst* **37**, 124575–124611 (2024)
- [19] Glicksberg, B.S., Timsina, P., Patel, D., Sawant, A., Vaid, A., Raut, G., Charney, A.W., Apakama, D., Carr, B.G., Freeman, R., Nadkarni, G.N., Klang, E.: Evaluating the accuracy of a state-of-the-art large language model for prediction of admissions from the emergency room. *J. Am. Med. Inform. Assoc.* **31**(9), 1921–1928 (2024)
- [20] Nazyrova, N., Chahed, S., Chausalet, T., Dwek, M.: Leveraging large language models for medical text classification: a hospital readmission prediction case. In: 2024 14th International Conference on Pattern Recognition Systems (ICPRS), pp. 1–7. IEEE, ??? (2024)
- [21] Ben Shoham, O., Rappoport, N.: Cpllm: Clinical prediction with large language models. *PLOS Digital Health* **3**(12), 0000680 (2024) <https://doi.org/10.1371/journal.pdig.0000680> . Published December 6, 2024
- [22] Scarlat, A., Campion, F.X.: Predicting 30-day mortality and readmission using hospital discharge summaries: A comparative analysis of machine learning models, large language models, and physicians. *medRxiv*, 2025–032625324714 (2025)
- [23] Bhasuran, B., Jin, Q., Xie, Y., Yang, C., Hanna, K., Costa, J., Shavor, C., Han, W., Lu, Z., He, Z.: Preliminary analysis of the impact of lab results on large language model generated differential diagnoses. *NPJ Digit. Med.* **8**(1), 166 (2025)
- [24] McDuff, D., Schaekermann, M., Tu, T., Palepu, A., Wang, A., Garrison, J., Singhal, K., Sharma, Y., Azizi, S., Kulkarni, K., Hou, L., Cheng, Y., Liu, Y., Mahdavi, S.S., Prakash, S., Pathak, A., Semturs, C., Patel, S., Webster, D.R., Dominowska, E., Gottweis, J., Barral, J., Chou, K., Corrado, G.S., Matias, Y., Sunshine, J., Karthikesalingam, A., Natarajan, V.: Towards accurate differential diagnosis with large language models. *Nature*, 1–7 (2025)
- [25] Alyakin, A., Stryker, J., Alber, D.A., Sangwon, K.L., Lee, J.V., Duderstadt, B., Save, A., Kurland, D., Frome, S., Singh, S., Zhang, J., Yang, E., Park, K.Y., Orillac, C., Valliani, A.A., Neifert, S., Liu, A., Patel, A., Livia, C., Lau, D., Laufer, I., Rozman, P.A., Hidalgo, E.T., Riina, H., Feng, R., Hollon, T., Aphinyanaphongs, Y., Golfinos, J.G., Snyder, L., Leuthardt, E., Kondziolka, D., Oermann, E.K.: CNS-Obsidian: A Neurosurgical Vision-Language Model Built From Scientific Publications (2025). <https://arxiv.org/abs/2502.19546>
- [26] Nori, H., Lee, Y.T., Zhang, S., Carignan, D., Edgar, R., Fusi, N., King, N., Larson, J., Li, Y., Liu, W., Luo, R., McKinney, S.M., Ness, R.O., Poon, H., Qin, T., Usuyama, N., White, C., Horvitz, E.: Can generalist foundation models outcompete special-purpose tuning? case study in medicine. *ArXiv abs/2311.16452* (2023)

- [27] Zhou, H.-Y., Adithan, S., Acosta, J.N., Topol, E.J., Rajpurkar, P.: A generalist learner for multifaceted medical image interpretation. *arXiv [cs.CV]* (2024) [cs.CV]
- [28] Lehman, E., Hernandez, E., Mahajan, D., Wulff, J., Smith, M.J., Ziegler, Z., Nadler, D., Szolovits, P., Johnson, A., Alsentzer, E.: Do we still need clinical language models? *arXiv [cs.CL]* (2023) [cs.CL]
- [29] Johri, S., Jeong, J., Tran, B.A., Schlessinger, D.I., Wongvibulsin, S., Barnes, L.A., Zhou, H.-Y., Cai, Z.R., Van Allen, E.M., Kim, D., Daneshjou, R., Rajpurkar, P.: An evaluation framework for clinical use of large language models in patient interaction tasks. *Nature medicine* **31**(1), 77–86 (2025)
- [30] Jiang, Y., Black, K.C., Geng, G., Park, D., Zou, J., Ng, A.Y., Chen, J.H.: MedAgentBench: A realistic virtual EHR environment to benchmark medical LLM agents. *arXiv [cs.LG]* (2025) [cs.LG]
- [31] Bean, A.M., Payne, R., Parsons, G., Kirk, H.R., Ciro, J., Mosquera, R., Monsalve, S.H., Ekanayaka, A.S., Tarassenko, L., Rocher, L., Mahdi, A.: Clinical knowledge in llms does not translate to human interactions. *ArXiv* **abs/2504.18919** (2025)
- [32] Vishwanath, K., Alyakin, A., Alber, D.A., Lee, J.V., Kondziolka, D., Oermann, E.K.: Medical large language models are easily distracted (2025). <https://arxiv.org/abs/2504.01201>
- [33] Yan, C., Fu, X., Xiong, Y., Wang, T., Hui, S.C., Wu, J., Liu, X.: Llm sensitivity evaluation framework for clinical diagnosis. *ArXiv* **abs/2504.13475** (2025)
- [34] Hager, P., Jungmann, F., Holland, R., Bhagat, K., Hubrecht, I., Knauer, M., Vielhauer, J., Makowski, M., Braren, R., Kaissis, G., Rueckert, D.: Evaluation and mitigation of the limitations of large language models in clinical decision-making. *Nat. Med.* **30**(9), 2613–2622 (2024)
- [35] Arora, R.K., Wei, J., Hicks, R.S., Bowman, P., Quiñóner-Candela, J., Tsimpourlas, F., Sharman, M., Shah, M., Vallone, A., Beutel, A., Heidecke, J., Singhal, K.: HealthBench: Evaluating large language models towards improved human health. *arXiv [cs.CL]* (2025) [cs.CL]
- [36] Wornow, M., Xu, Y., Thapa, R., Patel, B., Steinberg, E., Fleming, S., Pfeffer, M.A., Fries, J., Shah, N.H.: The shaky foundations of large language models and foundation models for electronic health records. *NPJ Digit. Med.* **6**(1), 135 (2023)
- [37] Johnson, A.E.W., Bulgarelli, L., Shen, L., Gayles, A., Shammout, A., Horng, S., Pollard, T.J., Hao, S., Moody, B., Gow, B., Lehman, L.-w.H., Celi, L.A., Mark, R.G.: MIMIC-IV, a freely accessible electronic health record dataset. *Scientific Data* **10**(1) (2023) <https://doi.org/10.1038/s41597-022-01899-x>
- [38] pubmed.gov: Download PubMed Data. NCBI Literature Resources

- [39] Sushil, M., Ludwig, D., Butte, A.J., Rudrapatna, V.A.: Developing a general-purpose clinical language inference model from a large corpus of clinical notes. *arXiv [cs.CL]* (2022) [cs.CL]
- [40] Yang, X., Chen, A., PourNejatian, N., Shin, H.C., Smith, K.E., Parisien, C., Compas, C., Martin, C., Costa, A.B., Flores, M.G., Zhang, Y., Magoc, T., Harle, C.A., Lipori, G., Mitchell, D.A., Hogan, W.R., Shenkman, E.A., Bian, J., Wu, Y.: A large language model for electronic health records. *NPJ Digit. Med.* **5**(1), 194 (2022)
- [41] Peng, C., Yang, X., Chen, A., Smith, K.E., PourNejatian, N., Costa, A.B., Martin, C., Flores, M.G., Zhang, Y., Magoc, T., Lipori, G., Mitchell, D.A., Ospina, N.S., Ahmed, M.M., Hogan, W.R., Shenkman, E.A., Guo, Y., Bian, J., Wu, Y.: A study of generative large language model for medical research and healthcare. *NPJ Digit. Med.* **6**(1), 210 (2023)
- [42] Jiang, L.Y., Liu, X.C., Nejatian, N.P., Nasir-Moin, M., Wang, D., Abidin, A., Eaton, K., Riina, H.A., Laufer, I., Punjabi, P., Miceli, M., Kim, N.C., Orillac, C., Schnurman, Z., Livia, C., Weiss, H., Kurland, D., Neifert, S., Dastagirzada, Y., Kondziolka, D., Cheung, A.T.M., Yang, G., Cao, M., Flores, M., Costa, A.B., Aphinyanaphongs, Y., Cho, K., Oermann, E.K.: Health system-scale language models are all-purpose prediction engines. *Nature* (2023)
- [43] Goh, E., Gallo, R., Hom, J., Strong, E., Weng, Y., Kerman, H., Cool, J.A., Kanjee, Z., Parsons, A.S., Ahuja, N., Horvitz, E., Yang, D., Milstein, A., Olson, A.P.J., Rodman, A., Chen, J.H.: Large language model influence on diagnostic reasoning: A randomized clinical trial: A randomized clinical trial. *JAMA Netw. Open* **7**(10), 2440969 (2024)
- [44] Nori, H., Daswani, M., Kelly, C., Lundberg, S., Ribeiro, M.T., Wilson, M., Liu, X., Sounderajah, V., Carlson, J., Lungren, M.P., Gross, B., Hames, P., Suleyman, M., King, D., Horvitz, E.: Sequential diagnosis with language models. *arXiv [cs.CL]* (2025) [cs.CL]
- [45] Tu, T., Schaekermann, M., Palepu, A., Saab, K., Freyberg, J., Tanno, R., Wang, A., Li, B., Amin, M., Cheng, Y., Vedadi, E., Tomasev, N., Azizi, S., Singhal, K., Hou, L., Webson, A., Kulkarni, K., Mahdavi, S.S., Semturs, C., Gottweis, J., Barral, J., Chou, K., Corrado, G.S., Matias, Y., Karthikesalingam, A., Natarajan, V.: Towards conversational diagnostic artificial intelligence. *Nature*, 1–9 (2025)
- [46] Luccioni, A.S., Jernite, Y., Strubell, E.: Power hungry processing: Watts driving the cost of AI deployment? *arXiv [cs.LG]* (2023) [cs.LG]
- [47] Cottier, B., Rahman, R., Fattorini, L., Maslej, N., Besiroglu, T., Owen, D.: The rising costs of training frontier AI models. *arXiv [cs.CY]* (2024) [cs.CY]
- [48] Feng, S., Ding, W., Liu, A., Wang, Z., Shi, W., Wang, Y., Shen, Z., Han, X., Lang, H., Lee, C.-Y., Pfister, T., Choi, Y., Tsvetkov, Y.: When one LLM drools, multi-LLM collaboration rules. *arXiv [cs.CL]* (2025) [cs.CL]

- [49] Belcak, P., Heinrich, G., Diao, S., Fu, Y., Dong, X., Muralidharan, S., Lin, Y.C., Molchanov, P.: Small language models are the future of agentic AI. arXiv [cs.AI] (2025) [cs.AI]
- [50] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C.L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Aspell, A., Welinder, P., Christiano, P., Leike, J., Lowe, R.: Training language models to follow instructions with human feedback. arXiv [cs.CL] (2022) [cs.CL]
- [51] Rocklin, M.: Dask: Parallel computation with blocked algorithms and task scheduling. In: Proceedings of the Python in Science Conference, pp. 126–132. SciPy, ??? (2015)
- [52] Soboleva, D., Al-Khateeb, F., Myers, R., Steeves, J.R., Hestness, J., Dey, N.: SlimPajama: A 627B token cleaned and deduplicated version of RedPajama. <https://www.cerebras.net/blog/slimpajama-a-627b-token-cleaned-and-deduplicated-version-of-redpajama> (2023). <https://huggingface.co/datasets/cerebras/SlimPajama-627B>
- [53] Common Crawl Foundation: Common Crawl Dataset. <https://commoncrawl.org>. Accessed: 2025-11-14 (2024)
- [54] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J.: C4: Colossal Clean Crawled Corpus. <https://www.tensorflow.org/datasets/catalog/c4>. Accessed: YYYY-MM-DD (2020)
- [55] Wenzek, G., Lacroix, T., Lavergne, T., et al.: BookCorpus2. https://github.com/facebookresearch/cc_net. Included in The Pile and SlimPajama datasets. (2020)
- [56] Rae, J.W., Potapenko, A., Jayakumar, S.M., Lillicrap, T.: Compressing large-scale language models. In: International Conference on Machine Learning (ICML) (2020). PG-19 long-book subset from Project Gutenberg.
- [57] Charlson, M.E., Pompei, P., Ales, K.L., MacKenzie, C.R.: A new method of classifying prognostic comorbidity in longitudinal studies: Development and validation. *Journal of Chronic Diseases* **40**(5), 373–383 (1987) [https://doi.org/10.1016/0021-9681\(87\)90171-8](https://doi.org/10.1016/0021-9681(87)90171-8)
- [58] Charlson Comorbidity Index (CCI). <https://www.mdcalc.com/calc/3917/charlson-comorbidity-index-cci>. Accessed: 2025-10-12
- [59] Johnson, A., Pollard, T., Mark, R.: MIMIC-III clinical database. PhysioNet (2023)
- [60] Yang, G., Cao, M., Jiang, L.Y., Liu, X.C., Cheung, A.T.M., Weiss, H., Kurland, D., Cho, K., Oermann, E.K.: Language model classifier aligns better with physician word sensitivity than XGBoost on readmission prediction. arXiv [cs.CL] (2022) [cs.CL]
- [61] Welbl, J., Liu, N.F., Gardner, M.: Crowdsourcing multiple choice science questions. arXiv [cs.HC] (2017) [cs.HC]

- [62] Jin, Q., Dhingra, B., Liu, Z., Cohen, W.W., Lu, X.: PubMedQA: A dataset for biomedical research question answering. arXiv [cs.CL] (2019) [cs.CL]
- [63] Lightning AI: LitGPT. <https://github.com/Lightning-AI/litgpt> (2023)
- [64] Zhao, Y., Gu, A., Varma, R., Luo, L., Huang, C.-C., Xu, M., Wright, L., Shojanazeri, H., Ott, M., Shleifer, S., Desmaison, A., Balioglu, C., Damania, P., Nguyen, B., Chauhan, G., Hao, Y., Mathews, A., Li, S.: PyTorch FSDP: Experiences on scaling fully sharded data parallel. arXiv [cs.DC] (2023) [cs.DC]
- [65] Bengio, Y.: Practical recommendations for gradient-based training of deep architectures. arXiv [cs.LG] (2012) [cs.LG]
- [66] Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: LoRA: Low-rank adaptation of large language models. arXiv [cs.CL] (2021) [cs.CL]
- [67] Biderman, S., Schoelkopf, H., Sutawika, L., Gao, L., Tow, J., Abbasi, B., Aji, A.F., Ammanamanchi, P.S., Black, S., Clive, J., DiPofi, A., Etxaniz, J., Fattori, B., Forde, J.Z., Foster, C., Hsu, J., Jaiswal, M., Lee, W.Y., Li, H., Lovering, C., Muennighoff, N., Pavlick, E., Phang, J., Skowron, A., Tan, S., Tang, X., Wang, K.A., Winata, G.I., Yvon, F., Zou, A.: Lessons from the trenches on reproducible evaluation of language models. arXiv [cs.CL] (2024) [cs.CL]
- [68] Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C.H., Gonzalez, J.E., Zhang, H., Stoica, I.: Efficient memory management for large language model serving with PagedAttention. arXiv [cs.LG] (2023) [cs.LG]
- [69] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Müller, A., Nothman, J., Louppe, G., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, A.: Scikit-learn: Machine learning in python. arXiv [cs.LG] (2012) [cs.LG]
- [70] Detlefsen, N., Borovec, J., Schock, J., Jha, A., Koker, T., Di Liello, L., Stancil, D., Quan, C., Grechkin, M., Falcon, W.: TorchMetrics - measuring reproducibility in PyTorch. J. Open Source Softw. 7(70), 4101 (2022)