

Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement
- A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided
Only common tests should be described solely by name; describe more complex techniques in the Methods section.
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
Give P values as exact values whenever suitable.
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection

SQL scripts executed via Cloudera Hue were used to extract clinical notes from the NYU Langone Health electronic health record (EHR). Raw CSV exports were processed with Python Dask (v2023.12) for distributed loading and concatenation. Text cleaning (punctuation standardization, non-ASCII removal, whitespace normalization, short-note filtering) was performed using Python standard-library regular expressions. Public pretraining data included SlimPajama (627B tokens; Cerebras, 2023), downloaded and prepared using the LitGPT data preparation scripts. MIMIC III data were obtained from PhysioNet (v1.4). SciQ and PubMedQA evaluation datasets were accessed via the LM Evaluation Harness (EleutherAI). OASST2 instruction-tuning data was obtained from Hugging Face. No commercial or proprietary data-collection software beyond the EHR system was used.

Data analysis

Tokenization: Lang1 models used the SentencePiece tokenizer from Llama-2-7B (32K vocabulary). All other models used their respective pretrained tokenizers.

Pretraining and finetuning. Models were trained with PyTorch (v2.1) using the LitGPT library (v0.4) and Lightning Fabric with Fully Sharded Data Parallel (FSDP). Training was conducted on NVIDIA 80 GB H100 GPUs (8–64 GPUs). Configuration management used Hydra (v1.3). Hyperparameter search for finetuning used Optuna (v3.4) with a random sampler (5 trials per task). LoRA finetuning of Llama-70B used the PEFT library.

Inference. Zero-shot and few-shot evaluation used the LM Evaluation Harness (EleutherAI, v0.4). DeepSeek models were served via vLLM. On-premises GPT-4o was accessed via Azure OpenAI Service.

Statistical analysis. AUROC was computed with scikit-learn (v1.3). One-Versus-Rest AUROC was used for multiclass tasks (LOS, CCI). Calibration curves were generated with scikit-learn (n=15 bins). Expected calibration error (ECE) was calculated with torchmetrics (n=15 bins). 95%

confidence intervals were obtained by quantile bootstrap resampling (1,000 iterations) using SciPy (v1.11). Figures were generated with Matplotlib.

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our [policy](#)

These data consist of the production medical records of NYU Langone and cannot be made publicly available. Researchers may obtain a limited de-identified dataset (or a test subset) from NYU Langone Health System by reasonable request and subject to local and national ethical approvals. We also used publicly available i2b2-2012 (<https://portal.dbmi.hms.harvard.edu/projects/n2c2-nlp/>) and MIMIC-III (<https://physionet.org/content/mimiciii/1.4/>) datasets.

Research involving human participants, their data, or biological material

Policy information about studies with [human participants or human data](#). See also policy information about [sex, gender \(identity/presentation\), and sexual orientation](#) and [race, ethnicity and racism](#).

Reporting on sex and gender	Sex was determined from structured EHR fields (recorded biological sex). Model performance was stratified by sex (Male/Female) and reported in Extended Data (Stratified Evaluation). Sex was not used as a model input feature. Gender identity was not assessed as it is not reliably captured in the EHR.
Reporting on race, ethnicity, or other socially relevant groupings	Race and ethnicity categories were drawn from structured EHR fields as self-reported by patients during registration. Model performance was stratified by race, ethnicity, borough, age and pediatric status, with results reported in Extended Data. Race, ethnicity, and other demographic variables were not used as model input features. Stratified analyses were performed to assess equity of model performance across demographic groups.
Population characteristics	The study population comprises inpatients at NYU Langone Health (2003–2024) and Beth Israel Deaconess Medical Center (MIMIC III, 2001–2012). The pretraining corpus (NYU Notes+) includes 11,689,342 patients and 180,487,092 notes. Finetuning and evaluation datasets range from 87,974 to 421,429 patients per task. Detailed dataset statistics including note counts, patient counts, and class ratios are provided in Supplementary and NYU Notes+ demographics.
Recruitment	No prospective recruitment was performed. All data were collected retrospectively from existing EHR systems as part of routine clinical care. MIMIC III is a publicly available de-identified critical care database.
Ethics oversight	This study was approved by the Institutional Review Board at NYU Langone Health (protocol s21-01189). MIMIC III is governed by its own data use agreement via PhysioNet.

Note that full information on the approval of the study protocol must also be provided in the manuscript.

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences Behavioural & social sciences Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see nature.com/documents/nr-reporting-summary-flat.pdf

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	No statistical methods were used to predetermine sample size. For pretraining, we used all available clinical notes from the NYU Langone Health EHR (NYU Notes+: 11,689,342 patients, 180,487,092 notes, 2003–2023). For finetuning and evaluation, dataset sizes were determined by the availability of labelled encounters: readmission (604,326 notes), mortality (566,748 notes), LOS (566,748 notes), CCI (443,915 notes), and insurance denial (97,837 notes). External validation used MIMIC III (5,658–52,725 examples per task). These sample sizes are consistent with or larger than comparable clinical NLP studies (e.g., NYUTron).
Data exclusions	Rehabilitation, dialysis, and palliative care notes were excluded from the readmission dataset to focus on acute readmission, consistent with prior work (NYUTron). Notes shorter than 10 words or containing only placeholder values (e.g., <NA>) were excluded during preprocessing. For MIMIC III mortality/LOS, notes written more than 120 hours after admission were excluded. For stratified evaluation, demographic subgroups with single-class representation due to small sample size were omitted. All exclusion criteria were established before analysis.
Replication	Finetuning hyperparameter searches used 5 independent Optuna trials per task with random sampling. Bootstrap confidence intervals (1,000 resamples) were computed for all reported AUROCs. Key findings were replicated across multiple model scales (100M, 1B, 7B) and validated

on an external dataset (MIMIC III, Beth Israel Deaconess Medical Center). The pretraining dynamics result (clinical prediction does not emerge from pretraining alone) was consistent across all five tasks and all model sizes.

Randomization Finetuning datasets were split temporally rather than randomly: train/validation/test from 2013 to May 2021 (8:1:1 random split within this window), with separate temporal test sets from June–December 2021 and 2024. This temporal split design approximates real-world deployment conditions.

Blinding Blinding was not relevant to this study. All evaluations were performed computationally using predefined metrics (AUROC) on held-out test sets with no subjective human assessment.

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

Methods

- n/a Involved in the study
- Antibodies
- Eukaryotic cell lines
- Palaeontology and archaeology
- Animals and other organisms
- Clinical data
- Dual use research of concern
- Plants

- n/a Involved in the study
- ChIP-seq
- Flow cytometry
- MRI-based neuroimaging

Antibodies

Antibodies used *Describe all antibodies used in the study; as applicable, provide supplier name, catalog number, clone name, and lot number.*

Validation *Describe the validation of each primary antibody for the species and application, noting any validation statements on the manufacturer's website, relevant citations, antibody profiles in online databases, or data provided in the manuscript.*

Eukaryotic cell lines

Policy information about [cell lines and Sex and Gender in Research](#)

Cell line source(s) *State the source of each cell line used and the sex of all primary cell lines and cells derived from human participants or vertebrate models.*

Authentication *Describe the authentication procedures for each cell line used OR declare that none of the cell lines used were authenticated.*

Mycoplasma contamination *Confirm that all cell lines tested negative for mycoplasma contamination OR describe the results of the testing for mycoplasma contamination OR declare that the cell lines were not tested for mycoplasma contamination.*

Commonly misidentified lines
(See [ICLAC](#) register) *Name any commonly misidentified cell lines used in the study and provide a rationale for their use.*

Palaeontology and Archaeology

Specimen provenance *Provide provenance information for specimens and describe permits that were obtained for the work (including the name of the issuing authority, the date of issue, and any identifying information). Permits should encompass collection and, where applicable, export.*

Specimen deposition *Indicate where the specimens have been deposited to permit free access by other researchers.*

Dating methods *If new dates are provided, describe how they were obtained (e.g. collection, storage, sample pretreatment and measurement), where they were obtained (i.e. lab name), the calibration program and the protocol for quality assurance OR state that no new dates are provided.*

Tick this box to confirm that the raw and calibrated dates are available in the paper or in Supplementary Information.

Ethics oversight *Identify the organization(s) that approved or provided guidance on the study protocol, OR state that no ethical approval or guidance was required and explain why not.*

Note that full information on the approval of the study protocol must also be provided in the manuscript.

Animals and other research organisms

Policy information about [studies involving animals](#); [ARRIVE guidelines](#) recommended for reporting animal research, and [Sex and Gender in Research](#)

Laboratory animals	<i>For laboratory animals, report species, strain and age OR state that the study did not involve laboratory animals.</i>
Wild animals	<i>Provide details on animals observed in or captured in the field; report species and age where possible. Describe how animals were caught and transported and what happened to captive animals after the study (if killed, explain why and describe method; if released, say where and when) OR state that the study did not involve wild animals.</i>
Reporting on sex	<i>Indicate if findings apply to only one sex; describe whether sex was considered in study design, methods used for assigning sex. Provide data disaggregated for sex where this information has been collected in the source data as appropriate; provide overall numbers in this Reporting Summary. Please state if this information has not been collected. Report sex-based analyses where performed, justify reasons for lack of sex-based analysis.</i>
Field-collected samples	<i>For laboratory work with field-collected samples, describe all relevant parameters such as housing, maintenance, temperature, photoperiod and end-of-experiment protocol OR state that the study did not involve samples collected from the field.</i>
Ethics oversight	<i>Identify the organization(s) that approved or provided guidance on the study protocol, OR state that no ethical approval or guidance was required and explain why not.</i>

Note that full information on the approval of the study protocol must also be provided in the manuscript.

Clinical data

Policy information about [clinical studies](#)

All manuscripts should comply with the ICMJE [guidelines for publication of clinical research](#) and a completed [CONSORT checklist](#) must be included with all submissions.

Clinical trial registration	<i>Provide the trial registration number from ClinicalTrials.gov or an equivalent agency.</i>
Study protocol	<i>Note where the full trial protocol can be accessed OR if not available, explain why.</i>
Data collection	<i>Describe the settings and locales of data collection, noting the time periods of recruitment and data collection.</i>
Outcomes	<i>Describe how you pre-defined primary and secondary outcome measures and how you assessed these measures.</i>

Dual use research of concern

Policy information about [dual use research of concern](#)

Hazards

Could the accidental, deliberate or reckless misuse of agents or technologies generated in the work, or the application of information presented in the manuscript, pose a threat to:

No	Yes	
<input type="checkbox"/>	<input type="checkbox"/>	Public health
<input type="checkbox"/>	<input type="checkbox"/>	National security
<input type="checkbox"/>	<input type="checkbox"/>	Crops and/or livestock
<input type="checkbox"/>	<input type="checkbox"/>	Ecosystems
<input type="checkbox"/>	<input type="checkbox"/>	Any other significant area

Experiments of concern

Does the work involve any of these experiments of concern:

- | No | Yes | |
|--------------------------|--------------------------|---|
| <input type="checkbox"/> | <input type="checkbox"/> | Demonstrate how to render a vaccine ineffective |
| <input type="checkbox"/> | <input type="checkbox"/> | Confer resistance to therapeutically useful antibiotics or antiviral agents |
| <input type="checkbox"/> | <input type="checkbox"/> | Enhance the virulence of a pathogen or render a nonpathogen virulent |
| <input type="checkbox"/> | <input type="checkbox"/> | Increase transmissibility of a pathogen |
| <input type="checkbox"/> | <input type="checkbox"/> | Alter the host range of a pathogen |
| <input type="checkbox"/> | <input type="checkbox"/> | Enable evasion of diagnostic/detection modalities |
| <input type="checkbox"/> | <input type="checkbox"/> | Enable the weaponization of a biological agent or toxin |
| <input type="checkbox"/> | <input type="checkbox"/> | Any other potentially harmful combination of experiments and agents |

Plants

- Seed stocks** *Report on the source of all seed stocks or other plant material used. If applicable, state the seed stock centre and catalogue number. If plant specimens were collected from the field, describe the collection location, date and sampling procedures.*
- Novel plant genotypes** *Describe the methods by which all novel plant genotypes were produced. This includes those generated by transgenic approaches, gene editing, chemical/radiation-based mutagenesis and hybridization. For transgenic lines, describe the transformation method, the number of independent lines analyzed and the generation upon which experiments were performed. For gene-edited lines, describe the editor used, the endogenous sequence targeted for editing, the targeting guide RNA sequence (if applicable) and how the editor was applied.*
- Authentication** *Describe any authentication procedures for each seed stock used or novel genotype generated. Describe any experiments used to assess the effect of a mutation and, where applicable, how potential secondary effects (e.g. second site T-DNA insertions, mosaicism, off-target gene editing) were examined.*

ChIP-seq

Data deposition

- Confirm that both raw and final processed data have been deposited in a public database such as [GEO](#).
- Confirm that you have deposited or provided access to graph files (e.g. BED files) for the called peaks.

Data access links *For "Initial submission" or "Revised version" documents, provide reviewer access links. For your "Final submission" document, May remain private before publication. provide a link to the deposited data.*

Files in database submission *Provide a list of all files available in the database submission.*

Genome browser session *Provide a link to an anonymized genome browser session for "Initial submission" and "Revised version" documents only, to (e.g. [UCSC](#)) enable peer review. Write "no longer applicable" for "Final submission" documents.*

Methodology

- Replicates** *Describe the experimental replicates, specifying number, type and replicate agreement.*
- Sequencing depth** *Describe the sequencing depth for each experiment, providing the total number of reads, uniquely mapped reads, length of reads and whether they were paired- or single-end.*
- Antibodies** *Describe the antibodies used for the ChIP-seq experiments; as applicable, provide supplier name, catalog number, clone name, and lot number.*
- Peak calling parameters** *Specify the command line program and parameters used for read mapping and peak calling, including the ChIP, control and index files used.*
- Data quality** *Describe the methods used to ensure data quality in full detail, including how many peaks are at FDR 5% and above 5-fold enrichment.*
- Software** *Describe the software used to collect and analyze the ChIP-seq data. For custom code that has been deposited into a community repository, provide accession details.*

Flow Cytometry

Plots

Confirm that:

- The axis labels state the marker and fluorochrome used (e.g. CD4-FITC).
- The axis scales are clearly visible. Include numbers along axes only for bottom left plot of group (a 'group' is an analysis of identical markers).
- All plots are contour plots with outliers or pseudocolor plots.
- A numerical value for number of cells or percentage (with statistics) is provided.

Methodology

Sample preparation

Describe the sample preparation, detailing the biological source of the cells and any tissue processing steps used.

Instrument

Identify the instrument used for data collection, specifying make and model number.

Software

Describe the software used to collect and analyze the flow cytometry data. For custom code that has been deposited into a community repository, provide accession details.

Cell population abundance

Describe the abundance of the relevant cell populations within post-sort fractions, providing details on the purity of the samples and how it was determined.

Gating strategy

Describe the gating strategy used for all relevant experiments, specifying the preliminary FSC/SSC gates of the starting cell population, indicating where boundaries between "positive" and "negative" staining cell populations are defined.

- Tick this box to confirm that a figure exemplifying the gating strategy is provided in the Supplementary Information.

Magnetic resonance imaging

Experimental design

Design type

Indicate task or resting state; event-related or block design.

Design specifications

Specify the number of blocks, trials or experimental units per session and/or subject, and specify the length of each trial or block (if trials are blocked) and interval between trials.

Behavioral performance measures

State number and/or type of variables recorded (e.g. correct button press, response time) and what statistics were used to establish that the subjects were performing the task as expected (e.g. mean, range, and/or standard deviation across subjects).

Acquisition

Imaging type(s)

Specify: functional, structural, diffusion, perfusion.

Field strength

Specify in Tesla

Sequence & imaging parameters

Specify the pulse sequence type (gradient echo, spin echo, etc.), imaging type (EPI, spiral, etc.), field of view, matrix size, slice thickness, orientation and TE/TR/flip angle.

Area of acquisition

State whether a whole brain scan was used OR define the area of acquisition, describing how the region was determined.

Diffusion MRI

Used

Not used

Preprocessing

Preprocessing software

Provide detail on software version and revision number and on specific parameters (model/functions, brain extraction, segmentation, smoothing kernel size, etc.).

Normalization

If data were normalized/standardized, describe the approach(es): specify linear or non-linear and define image types used for transformation OR indicate that data were not normalized and explain rationale for lack of normalization.

Normalization template

Describe the template used for normalization/transformation, specifying subject space or group standardized space (e.g. original Talairach, MNI305, ICBM152) OR indicate that the data were not normalized.

Noise and artifact removal

Describe your procedure(s) for artifact and structured noise removal, specifying motion parameters, tissue signals and physiological signals (heart rate, respiration).

Volume censoring

*Define your software and/or method and criteria for volume censoring, and state the extent of such censoring.***Statistical modeling & inference**

Model type and settings

Specify type (mass univariate, multivariate, RSA, predictive, etc.) and describe essential details of the model at the first and second levels (e.g. fixed, random or mixed effects; drift or auto-correlation).

Effect(s) tested

*Define precise effect in terms of the task or stimulus conditions instead of psychological concepts and indicate whether ANOVA or factorial designs were used.*Specify type of analysis: Whole brain ROI-based Both

Statistic type for inference

Specify voxel-wise or cluster-wise and report all relevant parameters for cluster-wise methods.(See [Eklund et al. 2016](#))

Correction

*Describe the type of correction and how it is obtained for multiple comparisons (e.g. FWE, FDR, permutation or Monte Carlo).***Models & analysis**

n/a | Involved in the study

 Functional and/or effective connectivity Graph analysis Multivariate modeling or predictive analysis

Functional and/or effective connectivity

Report the measures of dependence used and the model details (e.g. Pearson correlation, partial correlation, mutual information).

Graph analysis

Report the dependent variable and connectivity measure, specifying weighted graph or binarized graph, subject- or group-level, and the global and/or node summaries used (e.g. clustering coefficient, efficiency, etc.).

Multivariate modeling and predictive analysis

Specify independent variables, features extraction and dimension reduction, model, training and evaluation metrics.