

MEDS Decentralized, Extensible Validation (MEDS-DEV) Benchmark: Establishing Reproducibility and Comparability in ML for Health

This effort is a collaboration of many groups and individuals. In this listing, we list working groups within this effort, and individuals involved in those groups in alphabetical order by first name.

The Core MEDS-DEV Working Group:

Aleksia Kolo, Chao Pang, Edward Choi, Ethan Steinberg, Hyewon Jeong, Jack Gallifant, Jason A Fries, Jeffrey N. Chiang, Jungwoo Oh, Justin Xu, Kamilė Stankevičiūtė, Kiril V. Klein, Matthew B. A. McDermott, Mikkel Odgaard, Nassim Oufattole, Nigam H. Shah, Patrick Rockenschaub, Pawel Renc, Robin P. van de Water, Shalmali Joshi, Simon A. Lee, Teya S. Bergamaschi, Tom J. Pollard, Vincent Jeanselme, Young Sang Choi

MATTMCDERMOTT8@GMAIL.COM

Core MEDS-DEV Tools: ACES, FEMR, MEDS-Evaluation, MEDS-Reader, MEDS-Transforms

Chao Pang, Ethan Steinberg, Hyewon Jeong, Jack Gallifant, Jason A Fries, Justin Xu, Kamilė Stankevičiūtė, Matthew B. A. McDermott, Michael Wornow, Nassim Oufattole, Nigam H. Shah, Pawel Renc, Teya S. Bergamaschi, Vincent Jeanselme

MIMIC-IV MEDS Dataset

Ethan Steinberg, Matthew B. A. McDermott, Nassim Oufattole, Pawel Renc, Tom J. Pollard

Columbia MEDS Dataset

Apara Kashyap, Chao Pang, Shalmali Joshi, Vincent Jeanselme, Xinzhuo Jiang, Yanwei Li, Young Sang Choi, Yuta Kobayashi

Other MEDS-DEV Ready Datasets (AUMCdb & eICU)

Matthew B. A. McDermott, Patrick Rockenschaub, Robin P. van de Water, Ryan C King

Profiled MEDS Models

Aleksia Kolo, Chao Pang, Edward Choi, Ethan Steinberg, Hyewon Jeong, Jason A Fries, Jungwoo Oh, Michael Wornow, Nassim Oufattole, Teya S. Bergamaschi, Xinzhuo Jiang

1. Introduction

Standardized benchmarks have driven significant progress in machine learning (ML) (Donoho, 2024; Deng et al., 2009). Benchmarks establish clearly defined metrics for success, allow researchers to fairly compare methods, and facilitate reproducibility and open science. Despite this, benchmarking remains underdeveloped in ML for healthcare.

Benchmarking in healthcare is hard for several reasons: (a) the lack of standardized schemas for sharing and processing data, which prohibit frictionless reproductions of published models over private health datasets; (b) ambiguity in how tasks and labels are defined, leading to irreproducible task definitions across papers; and (c) the inability to meaningfully compare model performance across the fragmented health data landscape Johnson et al. (2017); McDermott et al. (2021); Wang et al. (2020); Liao and Voldman (2023); Harutyunyan et al. (2019).

To resolve these limitations, we propose the MEDS Decentralized, Extensible Validation (MEDS-DEV) benchmark, a distributed benchmarking framework that enables seamless reproduction of model results with conceptually identical task definitions across a diverse set of source datasets, including both public and private datasets. MEDS-DEV differs from traditional benchmarks in a number of ways in order to be best suited to the ML4H ecosystem:

Decentralized Evaluation By default in MEDS-DEV, data is *not* presumed to be shareable or publicly available, and as such model architectures¹ will be evaluated on different datasets in a sparse, decentralized fashion driven by local collaborations and, eventually, larger competitions and curated efforts.

Extensible Task Landscape Secondly, MEDS-DEV is designed to operate over a large number of community curated, clinically meaningful tasks that are consistently defined in a dataset-agnostic manner. This permits the benchmark to both expand to diverse clinical areas of interest and to cover a much more rigorously curated and refined set of tasks through community engagement.

Rigorous, Comparable Validation Finally, by virtue of the MEDS standard and the seamless reproducibility it offers, MEDS-DEV can use identical evaluation systems across submitted models and

1. Note that MEDS-DEV primarily helps compare model architectures and training recipes rather than pre-trained models, given the lack of widespread data availability.

tasks, streamlining analysis of not only performance, but fairness metrics, calibration, computational costs of training and evaluation, dataset-size sensitivity, and more. This offers a significantly expanded set of analysis opportunities for models in the ML4H field. While only a subset of these evaluation metrics are currently implemented, all are clearly operationalizable in the MEDS-DEV model and on the planned future roadmap.

2. Method

Adding new results For a new task, model, and dataset combination, the results can be submitted via a pull request to the [MEDS-DEV GitHub repository](#) in the form of a standardized [MEDS-Evaluation](#) results file. The pull request will be reviewed by the maintainers of MEDS-DEV and the results incorporated into the leaderboard upon approval.

Adding new tasks MEDS-DEV prediction tasks are defined through the ACES configuration system (Xu et al., 2024). This library enables the expression of task predicates (e.g. phenotype and event definitions) and inclusion/exclusion criteria in a dataset-agnostic way. To propose a new task in MEDS-DEV, users submit a pull request to the [GitHub repository](#) with an ACES configuration file and (optionally) any dataset-specific mappings of ACES predicates in the corresponding dataset configuration files. The pull request should also contain a README describing the task and its associated clinical utility. The community can comment on this pull request to suggest changes to the task or to contest its clinical utility before its final inclusion in MEDS-DEV. Note that adding a new task to MEDS-DEV does not inherently generate new model results for that task; instead, the results can be added over time, reflecting the sparse, decentralized nature of MEDS-DEV. We hope that these clear, standardized and reproducible task definitions is a valuable contribution to ML4H given its current challenges with reproducing basic concepts such as mortality (Johnson et al., 2017).

Adding new models For participation in MEDS-DEV, a model must be runnable on any MEDS dataset with its outputs conforming to the standardized [MEDS-Evaluation](#) schema. Model description is submitted via a pull request to the [GitHub repository](#) along with any relevant metadata and usage instructions as required by the pull request template.

Model	MIMIC-IV/Columbia	
	Long LOS	ICU Mortality
Log. Reg.	0.752/0.677	0.754/0.509
LightGBM	0.783/0.757	0.798/0.661
MOTOR	0.804/0.735	0.854/0.727
CEHR-BERT	0.808/0.741	0.845/0.726
MEDS-Tab	0.811/0.761	0.830/ 0.785
GenHPF	0.779/0.662	0.790/0.633

Table 1: Proof of viability results (AUC-ROC). Results format: “MIMIC-IV”/“Columbia”.

94 **Adding new datasets** To support a MEDS-
 95 compatible dataset in MEDS-DEV, the contributor
 96 must define a predicate configuration file mapping
 97 dataset-specific features to task-specific concepts as
 98 defined by their ACES configuration files, document
 99 any limitations and incompatibilities (in terms of cen-
 100 soring, inclusion/exclusion criteria, and other poten-
 101 tial biases), and describe the access policy. As with
 102 other contributions, this information is submitted as
 103 a pull request to the [MEDS-DEV GitHub repository](#).

104 3. Results

105 MEDS-DEV is designed for extensibility and com-
 106 munity contribution; however, we have a number
 107 of existing, proof-of-viability results demonstrating
 108 this style of benchmarking, with reproducibility as
 109 a first class citizen. In particular, MEDS-compliant
 110 datasets for use in MEDS-DEV have already been
 111 curated for a number of public and private datasets,
 112 including (*public*) MIMIC-IV ([Johnson et al., 2023](#)),
 113 eICU ([Pollard et al., 2018](#)), AUMCdb ([Thoral et al.,](#)
 114 [2021](#)), EHRShot ([Wornow et al., 2023](#)), (*private*)
 115 Stanford data, Columbia data, and cohort-specific
 116 datasets from Toronto, Copenhagen, and Mass Gen-
 117 eral Brigham, with further datasets still under con-
 118 struction. In addition, MEDS-complaint versions
 119 of various published model architectures, such as
 120 MOTOR ([Steinberg et al., 2023](#)), CLMBR ([Stein-](#)
 121 [berg et al., 2021](#)), EBCL ([Jeong et al., 2024](#)),
 122 GenHPF ([Hur et al., 2023](#)), CEHR-BERT ([Pang](#)
 123 [et al., 2021](#)), and MEDS-Tab ([Oufattole et al., 2024](#))
 124 already in use, with further models actively being
 125 converted for inclusion, such as ETHOS ([Renc et al.,](#)
 126 [2024](#)), ESGPT ([McDermott et al., 2023](#)), CORE-
 127 BEHRT ([Odgaard et al., 2024](#)) and more.

128 For all of these models and datasets, MEDS-DEV
 129 contains a preliminary collection of 12 tasks across
 130 different clinical areas and challenges. MEDS-DEV
 131 is designed for the set of examined tasks to grow and
 132 change over time through community contribution,
 133 and we have already seen some of the most exten-
 134 sive and involved community discussions on the mer-
 135 its of different task inclusion/exclusion criteria in the
 136 [GitHub Issues for MEDS-DEV](#) that most authors in
 137 this project have encountered professionally to date.

138 To demonstrate the viability of transporting these
 139 models across public and private datasets, we show
 140 a subset of preliminary model results from MEDS-
 141 DEV in Table 1. This table shows comparison of
 142 newly trained models across 6 different model archi-
 143 tectures from 4 different author groups across a
 144 public and private dataset (MIMIC-IV and a subset
 145 of Columbia data, respectively), demonstrating that
 146 these models can be reliably reproduced across sites
 147 via the MEDS and MEDS-DEV frameworks. Note
 148 that these results are preliminary—and in particular
 149 the Columbia data used is only a 10K patient sub-
 150 set of their entire cohort—but they nevertheless es-
 151 tablish the viability of the MEDS-DEV system, thus
 152 motivating its presentation as a demonstration to the
 153 ML4H community to help encourage this new, signifi-
 154 cantly more reproducible style of learning in our field.

155 4. Discussion

156 MEDS-DEV represents a first step towards building a
 157 standardized benchmarking infrastructure for health-
 158 care research. It addresses the three main limitations
 159 of prior benchmarks: data standardization, task def-
 160 inition consistency, and multi-institution participa-
 161 tion. As shown in Section 3, MEDS-DEV enabled
 162 us, for the first time, to quickly evaluate four state-of-
 163 the-art EHR foundation models across multiple insti-
 164 tutions on a common set of tasks. As health systems
 165 begin to deploy models into the clinic, benchmarking
 166 efforts such as MEDS-DEV will serve an increasingly
 167 important role in validating models and help accel-
 168 erate the development of ML methods for EHR data.

169 With several dozen members already in the MEDS-
 170 DEV community, we are excited to build upon this
 171 momentum and present MEDS-DEV to the broader
 172 ML4H audience. We invite anyone who resonates
 173 with our vision for more rigorous, reproducible sci-
 174 ence to join the MEDS-DEV community and con-
 175 tribute models, datasets, and tasks here: <https://github.com/mmcdermott/MEDS-DEV>.

References

- 177
- 178 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai
179 Li, and Li Fei-Fei. Imagenet: A large-scale hier-
180 archical image database. In *2009 IEEE conference*
181 *on computer vision and pattern recognition*, pages
182 248–255. Ieee, 2009.
- 183 David Donoho. Data science at the singularity. *Har-*
184 *vard Data Science Review*, 6(1), 2024.
- 185 Hrayr Harutyunyan, Hrant Khachatrian, David C
186 Kale, Greg Ver Steeg, and Aram Galstyan. Multi-
187 task learning and benchmarking with clinical time
188 series data. *Scientific data*, 6(1):96, 2019.
- 189 Kyunghoon Hur, Jungwoo Oh, Junu Kim, Jiyoun
190 Kim, Min Jae Lee, Eunbyeol Cho, Seong-Eun
191 Moon, Young-Hak Kim, Louis Atallah, and Ed-
192 ward Choi. Genhpf: General healthcare predic-
193 tive framework for multi-task multi-source learn-
194 ing. *IEEE Journal of Biomedical and Health In-*
195 *formatics*, 2023.
- 196 Hyewon Jeong, Nassim Oufattole, Matthew Mcder-
197 mott, Aparna Balagopalan, Bryan Jangeesingh,
198 Marzyeh Ghassemi, and Collin Stultz. Event-based
199 contrastive learning for medical time series, 2024.
200 URL <https://arxiv.org/abs/2312.10308>.
- 201 Alistair EW Johnson, Tom J Pollard, and Roger G
202 Mark. Reproducibility in critical care: a mortal-
203 ity prediction case study. In *Machine learning for*
204 *healthcare conference*, pages 361–376. PMLR, 2017.
- 205 Alistair EW Johnson, Lucas Bulgarelli, Lu Shen,
206 Alvin Gayles, Ayad Shammout, Steven Horng,
207 Tom J Pollard, Sicheng Hao, Benjamin Moody,
208 Brian Gow, et al. MIMIC-IV, a freely accessible elec-
209 tronic health record dataset. *Scientific data*, 10(1):
210 1, 2023.
- 211 Wei Liao and Joel Voldman. A multidatabase extrac-
212 tion pipeline (metre) for facile cross validation in
213 critical care research. *Journal of Biomedical Infor-*
214 *matics*, 141:104356, 2023.
- 215 Matthew McDermott, Bret Nestor, Evan Kim, Wan-
216 cong Zhang, Anna Goldenberg, Peter Szolovits,
217 and Marzyeh Ghassemi. A comprehensive ehr time-
218 series pre-training benchmark. In *Proceedings of*
219 *the Conference on Health, Inference, and Learning*,
220 pages 257–278, 2021.
- 221 Matthew McDermott, Bret Nestor, Peniel Ar-
222 gaw, and Isaac S Kohane. Event stream gpt:
223 A data pre-processing and modeling library
224 for generative, pre-trained transformers over
225 continuous-time sequences of complex events. In
226 A. Oh, T. Naumann, A. Globerson, K. Saenko,
227 M. Hardt, and S. Levine, editors, *Advances*
228 *in Neural Information Processing Systems*,
229 volume 36, pages 24322–24334. Curran Asso-
230 ciates, Inc., 2023. URL [https://proceedings.
neurips.cc/paper_files/paper/2023/file/
4c8f197b24e9b05d22028c2de16a45d2-Paper-Datasets_
and_Benchmarks.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/4c8f197b24e9b05d22028c2de16a45d2-Paper-Datasets_and_Benchmarks.pdf). 231 233
- 234 Mikkel Odgaard, Kiril Vadimovic Klein, Sanne Møller
235 Thysen, Espen Jimenez-Solem, Martin Sillesen,
236 and Mads Nielsen. Core-behrt: A carefully op-
237 timized and rigorously evaluated behrt. *arXiv*
238 *preprint arXiv:2404.15201*, 2024.
- 239 Nassim Oufattole, Teya Bergamaschi, Aleksia Kolo,
240 Hyewon Jeong, Hanna Gaggin, Collin Stultz, and
241 Matthew McDermott. Meds tab: Scalable tabu-
242 larization and tabular feature usage utilities over
243 generic meds datasets. [https://github.com/
mmcdermott/MEDS_Tabular_AutoML](https://github.com/mmcdermott/MEDS_Tabular_AutoML), 2024. Ac-
244 cessed: 2024-10-31. 245
- 246 Chao Pang, Xinzhuo Jiang, Krishna S Kalluri,
247 Matthew Spotnitz, RuiJun Chen, Adler Perotte,
248 and Karthik Natarajan. Cehr-bert: Incorporating
249 temporal information from structured ehr data to
250 improve prediction tasks. In *Machine Learning for*
251 *Health*, pages 239–260. PMLR, 2021.
- 252 Tom J Pollard, Alistair EW Johnson, Jesse D Raffa,
253 Leo A Celi, Roger G Mark, and Omar Badawi.
254 The eicu collaborative research database, a freely
255 available multi-center database for critical care re-
256 search. *Scientific data*, 5(1):1–13, 2018.
- 257 Pawel Renc, Yugang Jia, Anthony E Samir, Jaroslaw
258 Was, Quanzheng Li, David W Bates, and Arka-
259 dius Sitek. Zero shot health trajectory prediction
260 using transformer. *NPJ Digital Medicine*, 7(1):256,
261 2024.
- 262 Ethan Steinberg, Ken Jung, Jason A Fries, Conor K
263 Corbin, Stephen R Pfohl, and Nigam H Shah.
264 Language models are an effective representation
265 learning technique for electronic health record
266 data. *Journal of biomedical informatics*, 113:
267 103637, 2021.

- 268 Ethan Steinberg, Jason Fries, Yizhe Xu, and Nigam
269 Shah. Motor: A time-to-event foundation model
270 for structured medical records. *arXiv preprint*
271 *arXiv:2301.03150*, 2023.
- 272 Patrick J Thorat, Jan M Peppink, Ronald H Driessen,
273 Eric JG Sijbrands, Erwin JO Kompanje, Lewis Ka-
274 plan, Heatherlee Bailey, Jozef Kesecioglu, Maurizio
275 Cecconi, Matthew Churpek, et al. Sharing icu pa-
276 tient data responsibly under the society of critical
277 care medicine/european society of intensive care
278 medicine joint data science collaboration: the am-
279 sterdam university medical centers database (ams-
280 terdamumcdb) example. *Critical care medicine*, 49
281 (6):e563–e577, 2021.
- 282 Shirly Wang, Matthew BA McDermott, Geet-
283 icka Chauhan, Marzyeh Ghassemi, Michael C
284 Hughes, and Tristan Naumann. Mimic-extract: A
285 data extraction, preprocessing, and representation
286 pipeline for mimic-iii. In *Proceedings of the ACM*
287 *conference on health, inference, and learning*, pages
288 222–235, 2020.
- 289 Michael Wornow, Rahul Thapa, Ethan Steinberg, Ja-
290 son Fries, and Nigam Shah. Ehrshot: An ehr
291 benchmark for few-shot evaluation of foundation
292 models. 2023.
- 293 Justin Xu, Jack Gallifant, Alistair EW Johnson, and
294 Matthew McDermott. Aces: Automatic cohort ex-
295 traction system for event-stream datasets. *arXiv*
296 *preprint arXiv:2406.19653*, 2024.