
Data Sheet for CT-Repo Dataset

1 **Contact:** Nafis Neehal (neehan@rpi.edu), Kristin P. Bennett (bennek@rpi.edu)

2 1 Data Sheet

3 Motivation

4 1. **For what purpose was the dataset created? (Was there a specific task in mind? Was**
5 **there a specific gap that needed to be filled? Please provide a description.)**

6 The dataset was developed to evaluate the performance of Large Language Models (LLMs)
7 in predicting baseline features using clinical trial metadata. The dataset includes baseline
8 features collected through the `clinicaltrials.gov` API as target. Currently, no other
9 available datasets provide a curated list of baseline features.

10 2. **Who created this dataset (e.g., which team, research group), and on behalf of which**
11 **entity (e.g., company, institution, organization)?**

12 The CT-Repo dataset was created by Nafis Neehal. At the time of creation, Nafis was a PhD
13 Candidate at RPI (CS).

14 3. **Who funded the creation of the dataset? (If there is an associated grant, please provide**
15 **the name of the grantor and the grant name and number.)**

16 The dataset creation was supported by IBM Research and the Rensselaer Institute for Data
17 Exploration and Applications.

18 4. **Any other comments?**

19 N/A

20 Composition

21 1. **What do the instances that comprise the dataset represent (e.g., documents, photos,**
22 **people, countries)? (Are there multiple types of instances (e.g., movies, users, and**
23 **ratings; people and interactions between them; nodes and edges)? Please provide a**
24 **description.)**

25 Each instance in the dataset represents a clinical trial study. It includes various textual
26 information about the clinical trial, such as the title, brief summary, conditions, inter-
27 ventions, primary outcome, eligibility criteria, and baseline features collected from the
28 `clinicaltrials.gov` API (see paper for details).

29 2. **How many instances are there in total (of each type, if appropriate)?**

30 There are a total of 1690 instances, each about a single clinical trial.

31 3. **Does the dataset contain all possible instances or is it a sample (not necessarily random)**
32 **of instances from a larger set? (If the dataset is a sample, then what is the larger set?**
33 **Is the sample representative of the larger set (e.g., geographic coverage)? If so, please**
34 **describe how this representativeness was validated/verified. If it is not representative of**

the larger set, please describe why not (e.g., to cover a more diverse range of instances, because instances were withheld or unavailable).)

The dataset contains all possible instances available at the time of crawling.

4. **What data does each instance consist of? ("Raw" data (e.g., unprocessed text or images) or features? In either case, please provide a description.)**

Each instance contains unprocessed texts for all features. See Table 1 in paper.

5. **Is there a label or target associated with each instance? If so, please provide a description.**

For the CT-Repo dataset, the target is to predict the baseline features collected through API and saved in the BaselineMeasures column.

6. **Is any information missing from individual instances? (If so, please provide a description, explaining why this information is missing (e.g., because it was unavailable). This does not include intentionally removed information, but might include, e.g., redacted text.)**

N/A

7. **Are relationships between individual instances made explicit (e.g., users' movie ratings, social network links)? (If so, please describe how these relationships are made explicit.)**

Instances are unrelated, each instance is about a separate clinical trial.

8. **Are there recommended data splits (e.g., training, development/validation, testing)? (If so, please provide a description of these splits, explaining the rationale behind them.)**

There are no data splits as no training/development/validation/testing is involved in our study.

9. **Are there any errors, sources of noise, or redundancies in the dataset? (If so, please provide a description.)**

Data has been curated to the best of our ability. We believe there are no further errors (removed a few erroneous keywords as baseline features, such as - 'Continuous', see paper) or redundancies (removed a few duplicate trials).

10. **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets, other datasets)? (If it links to or relies on external resources, a) are there guarantees that they will exist, and remain constant, over time; b) are there official archival versions of the complete dataset (i.e., including the external resources as they existed at the time the dataset was created); c) are there any restrictions (e.g., licenses, fees) associated with any of the external resources that might apply to a future user? Please provide descriptions of all external resources and any restrictions associated with them, as well as links or other access points, as appropriate.)**

Dataset is self-contained.

11. **Does the dataset contain data that might be considered confidential (e.g., data that is protected by legal privilege or by doctor-patient confidentiality, data that includes the content of individuals' non-public communications)? (If so, please provide a description.)**

No. All raw data in the dataset is from public sources (i.e. data from `clinicaltrials.gov` and publications).

12. **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety? (If so, please describe why.)**

N/A

13. **Does the dataset relate to people? (If not, you may skip the remaining questions in this section.)**

N/A

- 83 14. Does the dataset identify any subpopulations (e.g., by age, gender)? (If so, please
84 describe how these subpopulations are identified and provide a description of their
85 respective distributions within the dataset.)
86 N/A
- 87 15. Is it possible to identify individuals (i.e., one or more natural persons), either directly
88 or indirectly (i.e., in combination with other data) from the dataset? (If so, please
89 describe how.)
90 N/A
- 91 16. Does the dataset contain data that might be considered sensitive in any way (e.g., data
92 that reveals racial or ethnic origins, sexual orientations, religious beliefs, political
93 opinions or union memberships, or locations; financial or health data; biometric or
94 genetic data; forms of government identification, such as social security numbers;
95 criminal history)? (If so, please provide a description.)
96 N/A
- 97 17. Any other comments?
98 N/A

99 Collection Process

- 100 1. How was the data associated with each instance acquired? (Was the data directly ob-
101 servable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or
102 indirectly inferred/derived from other data (e.g., part-of-speech tags, model-based
103 guesses for age or language)? If data was reported by subjects or indirectly in-
104 ferred/derived from other data, was the data validated/verified? If so, please describe
105 how.)
106 Clinical Trial MetaData was reported in `clinicaltrials.gov` accessible through API.
- 107 2. What mechanisms or procedures were used to collect the data (e.g., hardware apparatus
108 or sensor, manual human curation, software program, software API)? (How were these
109 mechanisms or procedures validated?)
110 Trial metadata were collected using publicly available API.
- 111 3. If the dataset is a sample from a larger set, what was the sampling strategy (e.g.,
112 deterministic, probabilistic with specific sampling probabilities)?
113 The 1690 instances in the CT-Repo dataset are all the data available.
- 114 4. Who was involved in the data collection process (e.g., students, crowdworkers, contrac-
115 tors) and how were they compensated (e.g., how much were crowdworkers paid)?
116 One of the co-authors (Nafis Neehal) collected all the data.
- 117 5. Over what timeframe was the data collected? (Does this timeframe match the creation
118 timeframe of the data associated with the instances (e.g., recent crawl of old news
119 articles)? If not, please describe the timeframe in which the data associated with the
120 instances was created.)
121 The data was collected/curated during March-April 2024. However, the clinical trials
122 themselves have varying start and end dates, spanning several months/years.
- 123 6. Were any ethical review processes conducted (e.g., by an institutional review board)?
124 (If so, please provide a description of these review processes, including the outcomes,
125 as well as a link or other access point to any supporting documentation.)
126 N/A
- 127 7. Does the dataset relate to people? (If not, you may skip the remaining questions in this
128 section.)
129 N/A

- 130 8. **Did you collect the data from the individuals in question directly, or obtain it via third**
 131 **parties or other sources (e.g., websites)?**
 132 N/A
- 133 9. **Were the individuals in question notified about the data collection? (If so, please**
 134 **describe (or show with screenshots or other information) how notice was provided, and**
 135 **provide a link or other access point to, or otherwise reproduce, the exact language of**
 136 **the notification itself.)**
 137 N/A
- 138 10. **Did the individuals in question consent to the collection and use of their data? (If**
 139 **so, please describe (or show with screenshots or other information) how consent was**
 140 **requested and provided, and provide a link or other access point to, or otherwise**
 141 **reproduce, the exact language to which the individuals consented.)**
 142 N/A
- 143 11. **If consent was obtained, were the consenting individuals provided with a mechanism**
 144 **to revoke their consent in the future or for certain uses? (If so, please provide a**
 145 **description, as well as a link or other access point to the mechanism (if appropriate).)**
 146 N/A
- 147 12. **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g.,**
 148 **a data protection impact analysis) been conducted? (If so, please provide a description**
 149 **of this analysis, including the outcomes, as well as a link or other access point to any**
 150 **supporting documentation.)**
 151 N/A
- 152 13. **Any other comments?**
 153 N/A

154 Preprocessing/Cleaning/Labeling

- 155 1. **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucket-**
 156 **ing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances,**
 157 **processing of missing values)? (If so, please provide a description. If not, you may skip**
 158 **the remainder of the questions in this section.)**
 159 Originally, we started with around 1800 trials. After thorough preprocessing steps, including
 160 removing duplicate trials and those with missing values, we were left with 1693 trials for
 161 our final study (CT-Repo dataset). From these 1693 trials, we used 3 trials as examples for
 162 few-shot setting, and the remaining 1690 trials were used for the benchmarking purpose.
- 163 2. **Was the "raw" data saved in addition to the preprocessed/cleaned/labeled data (e.g., to**
 164 **support unanticipated future uses)? (If so, please provide a link or other access point**
 165 **to the "raw" data.)**
 166 Yes. Available in the same GitHub repository.
- 167 3. **Is the software used to preprocess/clean/label the instances available? (If so, please**
 168 **provide a link or other access point.)**
 169 Yes. The code is available in the GitHub repository.
- 170 4. **Any other comments?**
 171 N/A

172 Uses

- 173 1. **Has the dataset been used for any tasks already? (If so, please provide a description.)**
 174 The dataset has been used to benchmark State-of-the-art LLM's performance in predicting
 175 baseline features using clinical trial Metadata. We present detailed description of our
 176 experiment and data-usage throughout the paper.

- 177 2. **Is there a repository that links to any or all papers or systems that use the dataset? (If**
 178 **so, please provide a link or other access point.)**
 179 N/A - We are the first to release and use this dataset.
- 180 3. **What (other) tasks could the dataset be used for?**
 181 The dataset can be used for various studies, including making decisions about selecting
 182 different clinical trial design factors.
- 183 4. **Is there anything about the composition of the dataset or the way it was collected and**
 184 **preprocessed/cleaned/labeled that might impact future uses? (For example, is there**
 185 **anything that a future user might need to know to avoid uses that could result in unfair**
 186 **treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other**
 187 **undesirable harms (e.g., financial harms, legal risks) If so, please provide a description.**
 188 **Is there anything a future user could do to mitigate these undesirable harms?)**
 189 N/A
- 190 5. **Are there tasks for which the dataset should not be used? (If so, please provide a**
 191 **description.)**
 192 N/A
- 193 6. **Any other comments?**
 194 N/A

195 Distribution

- 196 1. **Will the dataset be distributed to third parties outside of the entity (e.g., company,**
 197 **institution, organization) on behalf of which the dataset was created? (If so, please**
 198 **provide a description.)**
 199 Yes, the dataset is freely available and accessible.
- 200 2. **How will the dataset be distributed (e.g., tarball on website, API, GitHub)? (Does the**
 201 **dataset have a digital object identifier (DOI)?)**
 202 Dataset is free for download at https://github.com/nafis-neeal/CTBench_LLM.
- 203 3. **When will the dataset be distributed?**
 204 The dataset is distributed as of June 2024 in its first version.
- 205 4. **Will the dataset be distributed under a copyright or other intellectual property (IP)**
 206 **license, and/or under applicable terms of use (ToU)? (If so, please describe this license**
 207 **and/or ToU, and provide a link or other access point to, or otherwise reproduce, any**
 208 **relevant licensing terms or ToU, as well as any fees associated with these restrictions.)**
 209 The dataset is distributed under CC0 1.0 Universal license.
- 210 5. **Have any third parties imposed IP-based or other restrictions on the data associated**
 211 **with the instances? (If so, please describe these restrictions, and provide a link or other**
 212 **access point to, or otherwise reproduce, any relevant licensing terms, as well as any**
 213 **fees associated with these restrictions.)**
 214 Not to our knowledge.
- 215 6. **Do any export controls or other regulatory restrictions apply to the dataset or to**
 216 **individual instances? (If so, please describe these restrictions, and provide a link or**
 217 **other access point to any supporting documentation.)**
 218 Not to our knowledge.
- 219 7. **Any other comments?**
 220 N/A

221 Maintenance

- 222 1. **Who is supporting/hosting/maintaining the dataset?**
223 Nafis Neehal is maintaining the dataset on his GitHub. Any further changes would be
224 announced through the GitHub repo link.
- 225 2. **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**
226 E-mail addresses are at the top of this document.
- 227 3. **Is there an erratum? (If so, please provide a link or other access point.)**
228 Currently, no additional versions are planned. However, if errors are encountered, future
229 versions of the dataset may be released (and will be versioned). All updates will be provided
230 in the same GitHub location with proper announcements.
- 231 4. **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete**
232 **instances)? (If so, please describe how often, by whom, and how updates will be**
233 **communicated to users (e.g., mailing list, GitHub)?)**
234 Same as previous.
- 235 5. **If the dataset relates to people, are there applicable limits on the retention of the data**
236 **associated with the instances (e.g., were individuals in question told that their data**
237 **would be retained for a fixed period of time and then deleted)? (If so, please describe**
238 **these limits and explain how they will be enforced.)**
239 No.
- 240 6. **Will older versions of the dataset continue to be supported/hosted/maintained? (If so,**
241 **please describe how. If not, please describe how its obsolescence will be communicated**
242 **to users.)**
243 Yes; all data will be versioned.
- 244 7. **If others want to extend/augment/build on/contribute to the dataset, is there a mecha-**
245 **nism for them to do so? (If so, please provide a description. Will these contributions be**
246 **validated/verified? If so, please describe how. If not, why not? Is there a process for**
247 **communicating/distributing these contributions to other users? If so, please provide a**
248 **description.)**
249 Errors may be submitted by emailing the authors. Further extensive augmentations may be
250 accepted at the authors' discretion.
- 251 8. **Any other comments?**
252 N/A