

Rundown of the Medical Natural Language Processing Task based on Actual Documents

S.M. Navin Nayer Anik
s.m.navin.nayer.anik@g.bracu.ac.bd

Sabrina Tabassum
sabrinatorabassum98@gmail.com

Fahim Abrar
fahim.abrar1@g.bracu.ac.bd

Md Farhadul Islam
md.farhadul.islam@g.bracu.ac.bd

Md Sabbir Hossain
md.sabbir.hossain1@g.bracu.ac.bd

Annajiat Alim Rasel
annajiat@gmail.com

Abstract—Information science, which includes natural language processing (NLP), significantly relies on standardized datasets for the development and evaluation of algorithms and models. In the medical field, privacy protection is a major concern, and acquiring datasets containing actual clinical data is difficult due to ethical and legal considerations. In order to advance NLP research in the medical domain, it is necessary to develop standardized and trustworthy medical datasets. In response to this need, we created the Real-MedNLP dataset for multiple medical tasks. Real-MedNLP aims to provide a real clinical dataset composed of radiology and case reports, as opposed to pseudo datasets that are frequently derived from medical textbooks or fictitious clinical texts. Moreover, the dataset contains both fundamental tasks, such as named entity recognition, and implemented practical tasks. Our work describes the Real-MedNLP evaluation task and the systems that were submitted for evaluation. The methods employed in the study adhere to a common paradigm founded on fundamental language models such as BERT, with the intent of addressing resource issues commonly encountered in NLP research. The study's findings shed light on the practicability of these approaches and their potential for future medical NLP research. Creating standardized medical datasets, such as Real-MedNLP, is crucial for the development and evaluation of NLP models in the medical field. Real-MedNLP's clinical data, bilingualism, and emphasis on both fundamental and applied practical tasks make it a valuable resource for future medical NLP research.

Index Terms—Medical Natural Language Processing, Named Entity Recognition, Case Reports, Radiography Reports, Adverse Drug Event, Case Identification

I. INTRODUCTION

There is an increasing demand for natural language processing (NLP) methods in healthcare due to the increased digitization of medical information. As natural language processing (NLP) becomes more widely discussed in the field of computer science, a rising number of researchers are turning to similar methods in the field of medicine. However, there is still a dearth of anonymized medical text data in languages other than English.

To look into this issue, four medical natural language processing (MedNLP) projects were completed: MedNLP-1, MedNLP-2, MedNLPPdoc, and MedWeb. This made it challenging to compare these tasks' results directly to those of other similar tasks carried out in English.

The Real-MedNLP dataset was developed to overcome these restrictions. Case reports and radiological reports are included

in this collection of authentic clinical literature. In addition to authentic other language data, the collection also has English translations of the original reports. Because of this, Real-MedNLP is the first standard for studying medical NLP across many languages.

The long-term objective of Real-MedNLP is to encourage the research and development of NLP systems suitable for use in actual medical environments. Therefore, the task scheme was modified to better represent the difficulties of dealing with clinical material in a healthcare setting. Among these obstacles is making sure the data are really valuable to physicians and processing text in a timely way.

In conclusion, the construction of the Real-MedNLP dataset fills a significant need in the study of medical NLP. Researchers may create and test NLP models useful in real-world healthcare settings with the help of Real-MedNLP's real-world clinical text data in many languages.

II. DATASET

The databases Cerner Health Facts® (version 2017) (Cerner) and Truven Health MarketScan® (Truven) were mined for our cohorts. These two datasets played significant roles in the study. More than 600 hospitals and clinics in the United States use the Cerner database, which is a repository for de-identified electronic health records (EHR). It contains longitudinal data spanning the years 2000 to 2017, representing over 68 million patients who are each entirely unique. The Truven Health MarketScan® Research Database (version 2015) includes an anonymized patient-level claims dataset. This information was removed to protect patient confidentiality. This study includes more than 170 million patients and spans the years 2011 through 2015. The data was obtained from insurance claims submitted by commercial insurance companies, Medicare supplement insurance providers, and Medicaid. The pretraining cohort used for Med-BERT was obtained from Cerner and contains a total of 28 million patients. We assembled three phenotyped cohorts to evaluate the model's validity and determine whether it is accurate. The DHF-Cerner cohort, the PaCa-Cerner cohort, and the PaCa-Truven cohort were all derived from Cerner. PaCa-Truven was the only cohort with a Truven origin. The results of the descriptive analysis performed on these cohorts are

shown in Table 2; for more information, see the "Methods" section. Refer to the page defining "cohort" for any additional information that may be necessary.

III. RELATED WORKS

In this [1] line of study, a Korean medical corpus is being compiled, and trained language models are being used. Deep learning natural language processing (NLP)-based model of the Korean medical language was offered by the authors of this work. The foundation for developing the model was taken from the BERT medical field pre-training program, and the model itself is predicated on a more intricate representation of the Korean language. Accuracy scores of 0.147 and 0.148 were achieved when the pre-trained model was applied to the masked language model with next sentence prediction. The table that follows presents these findings in tabular form. The fact that there was a considerable rise of 0.258 in the accuracy of the prediction of the following phrase during the intrinsic assessment might be viewed as a signal of progress. In addition, there was an increase of 0.053 in the Pearson correlation that was seen in the evaluation of Korean medical named entity recognition, and there was an increase of 0.046 that was seen in the extrinsic evaluation of Korean medical semantic textual similarity data. Both of these results can be attributed to the evaluation of Korean medical named entity recognition and Korean medical semantic textual similarity data. In the examination of the extrinsic evaluation of Korean medical named entity recognition, both of these outcomes were found. The investigation of the extrinsic data led to the discovery of both of these discoveries; hence, the conclusions that were formed from that data are relevant in this context.

The motivation for the suggested approach in this [2] study for problems relating to natural language processing (NLP), vision, and language is the use of self-supervised pre training of Transformer-style structures. In order to obtain deeper semantic representations of medical pictures and words, our technique includes using Masked Vision-Language Modeling (MVLM) as the pretext job on a large medical picture and caption dataset. This was done so that we could learn more about medical conditions. You may learn more about these depictions by looking at the subtitles. When applied to the two VQA datasets for radiology photographs (VQA-Med 2019 and VQA-RAD), the suggested method achieves results that are superior to those achieved by the ensemble models of the best answers achieved in the past. In addition to that, this technique provides attention maps that make the model interpretation process easier.

Using self-supervised pre-training of Transformer-style structures as its inspiration, this[3] study proposes a solution to the problems of NLP, vision, and language. Vision and Language Behind a Mask In order to accomplish the aim of developing richer semantic representations for medical images and writing using the methodology presented in this article, modelling is employed as the pretext job on a big dataset that is comprised of medical photographs and language. The performance of the suggested method is superior to that of

the ensemble models of the previous best solutions when it is applied to the twin VQA datasets for radiology photos, which are VQA-Med 2019 and VQA-RAD. In addition to that, this technique provides attention maps that make the model interpretation process easier.

In this[4] paper, a comprehensive literature review of work that makes use of deep learning for medical imaging and medical NLP is offered. The review covers a variety of topics, including tasks, pipelines, and impediments. The authors of this book have provided a complete examination of the architecture of deep learning as it is used in the disciplines of medical natural language processing and medical imaging. This analysis is offered in the book. The goal of this study is to improve diagnosis accuracy by determining the ideal combination of deep learning, natural language processing, and medical imaging. This work adds to the finding of the optimal combination. This research shed light on the considerable challenges that arise when using deep learning to medical natural language processing and medical imaging. These challenges have been brought to light as a result of the study.

The[5] research makes an effort to use NLP to forecast the medical specializations upon hospital admission in advance. The Amiens-Picardy University Hospital in France contributed more than 260K ED records to the study's massive dataset. The paper method seeks to combine structured data with free-form text notes made during the triage stage. On the one hand, the normal set of characteristics are tested using a conventional MLP model. On the other side, the textual data is processed using a convolutional neural network. Although each learning component is carried out separately and concurrently. The empirical findings showed generally accurate predictions. The study is thought to be another contribution to the growing efforts to use NLP techniques in the healthcare industry.

In the[6] research, an artificial intelligence-based medical chatbot was offered as a potential solution to this issue. This chatbot would be able to identify any illness and provide any relevant information about any disease. This chatbot's goals are to reduce healthcare costs and make medical information more approachable. This chatbot will serve as a virtual doctor, assisting patients in both diagnosing their illnesses and regaining their health. Only when a chatbot can accurately diagnose sickness and deliver the required information about the condition will a patient actually benefit from it. People may discuss their health issues with a text-to-text verdict bot, which also offers a customized diagnosis based on symptoms. People will therefore receive information about their health status and the appropriate level of safety. The writers of this article conducted a thorough review of recent literature. We looked at a number of papers over the past five years that are concerned with chatbots. A hybrid architecture built on deep learning models like NLP and the TF-IDF algorithm was also given in the article.

In this[7] research, a complete repository of medical abbreviations known as the Medical Abbreviation and Acronym Meta-Inventory was described. Eight source inventories from

various healthcare specializations and contexts were systematically combined to provide 104,057 abbreviations and 170,426 associated sensations. The application development process was sped up and redundancies were reduced thanks to cutting-edge machine learning technology, which enabled automatic cross-mapping of synonymous information. One of the added features is a quality check that is semi-automated, with the goal of reducing or eliminating errors. According to the Meta-Inventory, the completeness or coverage of abbreviations and senses in new clinical writing was far better than the next biggest repository. Specifically, the Meta-Inventory found that this was the case. This improvement was from 6% to 14 % in terms of abbreviation coverage, whereas it was between 28% to 52 % in terms of sense coverage. The Meta-Inventory is the most extensive and thorough compilation of medical abbreviations and acronyms written in dated American English to date. It is also the biggest collection of its kind. Applications may be found in a wide variety of contexts and areas as a result of the thorough coverage and diverse sources. The processing of natural language across institutions is now conceivable, while in the past this was not possible with earlier inventories.

The authors of this[8] paper perform a thorough analysis of several methods for incorporating medical knowledge into a pre-trained BERT model for clinical connection extraction. For the benchmark i2b2/VA 2010 clinical connection extraction dataset, the best model developed by the authors beats cutting-edge methods.

The authors of this[9] paper introduce Paper Plain1, a novel interactive interface that uses natural language processing to power four features: definitions of unfamiliar terms, in-situ plain language section summaries, a set of key questions that direct readers to the passages that contain the answers, and plain language summaries of the passages that contain the answers. Researchers who use Paper Plain have a better time reading and comprehending research papers without experiencing a reduction in paper comprehension compared to those who use a conventional PDF viewer. Overall, the findings of the study indicate that pointing readers to pertinent portions and offering "gists," or simple English summaries, alongside the actual paper material might facilitate reading medical papers and give readers more assurance when approaching them.

This[10] article includes information about the Medical Concept Annotation Analysis tool, which is available for free and open source (MedCAT). It includes (a) an innovative self-supervised machine learning algorithm for concept extraction utilizing any concept vocabulary, including UMLS/SNOMED-CT, (b) a feature-rich annotation interface for customizing and refining IE models, and (c) interfaces to the wider CogStack ecosystem for vendor-neutral health system implementation. In addition to all of these, we also provide an annotation interface that is loaded with features for the purpose of customization (c). It has been shown that ideas from the UMLS may now be obtained from open datasets. [Citation needed] [Citation needed] (F1:0.448–0.738 vs 0.429–0.650). SNOMED-CT extraction has been demonstrated at three of

London's most prestigious hospitals after first being validated in the real world through the use of self-supervised training on 8.8 billion words taken from 17 million clinical records. This was followed by additional fine-tuning through the utilization of 6,000 cases annotated by clinicians. By demonstrating a high degree of transferability (F1 \geq 0.94 across hospitals, datasets, and concept categories), the authors emphasize cross-domain EHR-agnostic value for clinical and research use cases. This value may be used for both clinical and research use cases.

IV. METHODOLOGY

A. Med-BERT architecture

The transformer design used in this study was almost identical to the one used in the aforementioned BERT paper. This design used both multilevel embeddings as well as bidirectional transformers for its implementation. We also used pre training procedures that were quite comparable (we used the same loss function for masking and classification before training). Because of conceptual distinctions among electronic health records (EHR) and text, it is not simple to apply the BERT methodology to structured EHR. In the first version of BERT, for example, the input medium was a 1-D chain of words, but ours is organized EHR, which is logged in a way that has multiple layers and relationships. Converting a structured EHR into a 1-D sequence or encoding its "structures" using the BERT transformer architecture is not well-defined. Both of these processes involve making the organized EHR easier to use. it understood what the necessary domain-specific pretraining operations are.

In order to fit the newly established modality, we presented our vision of the embedding layers for the Med-BERT system. To be more specific, the Med-BERT method used three different kinds of embeddings as inputs. We refer to the projections made from diagnostic codes as code embeddings, from the sequence of codes inside a visit as serialization embeddings, and from the location of a visit as visit embeddings. There are three types of embeddings: code embeddings, serialization embeddings, and visit embeddings. Visit embeddings identify each visit in the series, whereas code embeddings are low-dimensional representations of each diagnostic code. Serialization embeddings reflect the relative order or priority order of each code in each visit.

In contrast to BERT, we did not use CLS and SEP as input tokens in our algorithm. Our decision was primarily influenced by the disparity between EHR and text input formats. In BERT, each input sample contains only two adjacent sentences, and the token [SEP] serves as a separator between the two sentences for the pretraining task of next sentence prediction. This is done so the system can learn to forecast the following sentence. However, none of the assignments required us to predict the next sentence (this will be clarified in the following paragraph).

Based on our analysis, we determined that [SEP] is unnecessary because visit embeddings already sufficiently differentiate between visits. The main function of the [CLS] token in BERT

was to provide a concise overview of the material covered in the two sentences. In contrast, EHR sequences are usually significantly longer; for example, a series may comprise ten or more visits, and employing a single summary token would invariably result in the loss of data.

B. Pre Training Med-BERT

During the pre-training phase of our Med-BERT model, we made use of the same optimization approach and proposed hyperparameters (for more information, see "Implementation details") as were used during the pre-training phase of the original BERT model²⁹. In particular, this work was utilized to make an educated guess on the presence of BERT. More specifically, there was a possibility that a code would be replaced with [MASK 80A prognosis that the patient will have to remain in the hospital for a longer period of time than expected. We took the choice that for the classification task, rather than applying the question-answer pairs as in BERT, we would choose a clinical problem that had a relatively high frequency in our pre-training dataset and one that was not disease-specific. This decision was taken because we wanted to avoid the use of BERT's question-answer pairings. It was necessary to do these steps in order to ensure that our pre-trained model would have a greater degree of generalizability in the actual world. We identified and analyzed the three quality-of-care indicators that are used the most often; they are mortality, early readmission, and extended duration of stay in the hospital. Through comparison, we found that the objectives of mortality and early readmission are rather straightforward: the model quickly converges to an accuracy of more than 99%. As a direct consequence of this, we came to the conclusion that the prolonged LOS job would provide an excellent pre training opportunity. This assignment consisted of reviewing each patient to identify whether or not an event of a protracted hospital stay (LOS more than 7 days) had ever occurred during the course of the patient's EHR sequence. The LOS was calculated as the number of days the patient spent in the hospital divided by seven. We used this simplified model of extended LOS prediction by focusing at the patient level as opposed to the visit level in order to reduce the amount of complexity that was involved in the process of pretraining. The genuine future forecasting task is not something that we want to develop during the pre-training phase, and the masked LM task is also not something that we aim to perform. Neither of these things are anything that we want to do. The prolonged LOS job that is used for pretraining makes use of the bidirectional structure that is offered by Med-BERT, since this was the conclusion that we arrived at. A length of stay (LOS) that is longer than is typical not only reflects the patient's health status as it was recorded in earlier visits, but it also has an influence on visits that follow the one in which it occurs. On the other hand, activities such as disease start prediction or mortality will always be finished during the patient's most recent visit in the patient sequence. This is because these predictions are based on the most up-to-date

information available. This is due to the fact that the patient sequence's input data can only be constructed in a single way.

C. Refining Med-BERT for use in downstream prediction tasks

The pretraining fine-tuning paradigm is one that is adhered to by Med-BERT in a way that is comparable to that which is followed by BERT. The pretrained model will not construct anything other than a contextualized embedding for each and all of the input tokens that it is given. This is due to the fact that it has been instructed to behave in this manner. The model just generates an embedding, which may be used in a variety of applications; however, it does not directly give any prediction labels. For any downstream prediction task that is especially important, it is necessary to add a classification layer on top of the Med-BERT model. This is the need. A layer of this kind is sometimes referred to as a prediction head in certain circles. In order to make accurate predictions about the future, prediction heads are used. It is possible to add an extra layer of processing in the form of a simple prediction head, such as FFL, on top of the sequential output that is generated by the last layer of the Med-BERT algorithm. This option is available to users. A common kind of prediction head that is used for EHR predictive models is a recurrent neural network (RNN) that rolls over the output of token embeddings. Another name for this form of prediction head is a recurrent neural network (RNN). During the process of fine-tuning, we held loyal to the model that had been established by the first BERT, and we added a prediction head to the structure of the Med-BERT in order to expand its capabilities. In other words, we didn't deviate from the model that had been developed by the first BERT. After loading and initializing the parameters of the Med-BERT section by applying the pretrained model, the parameters of both the Med-BERT part and the prediction head were then updated by utilizing gradient descent. This took place after the loading and initialization of the Med-BERT section's parameters. This step was carried out after the previous one, which included the setup of the parameters for the Med-BERT component, was finished. This procedure was carried for an indefinite number of times until each of the parameters achieved their optimal state. The data that were used to develop the model originated from a training cohort that was particular to the condition in question; we referred to this cohort as the fine-tuning cohort. This group served as the primary contributor to the body of knowledge. This was a piece of information that was brought into the modeling process in the role of an input. The design of the model that had been pre trained and the model that had not been trained were identical. This design had randomly initiated embedding layers for tokens, segments, and positions, in addition to multi-head transformer tiers at various levels. We compared the results of the fine-tuning of the pretrained model to the results of the untrained model. This was done so that we could acquire knowledge of the extra values offered by the pretrained Med-BERT, most notably the utility of vast amounts of training data. This was done so that we could get knowledge of the additional values supplied by the pretrained Med-BERT. The

purpose for this was so that we could get this intelligence, and that was the driving force behind everything. This was done so that we could acquire a grasp of the extra values supplied by the pretrained Med-BERT, which was the purpose for doing this. This was done so that we could obtain a grasp of the additional values offered by the pretrained Med-BERT. The results of the performance of the models on the test set were used to fine-tune all of the models before going on to the test set, and the values that were presented were those results. The values that were shown were an illustration of the results that were obtained from this performance.

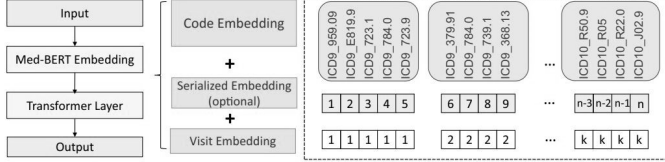


Fig. 1. title

D. Med-BERT's Visualization of Attentional Patterns

By using Med-BERT, it is possible to increase the accuracy of forecasts, in addition to making it simpler for people to understand such predictions. It is essential to investigate how the pretrained model has learnt by making use of the intricate structure and a significant amount of data. Doing so is not only intriguing but also helpful, thus it is essential that this be done. We provide a range of potential instances of how codes might be matched with one another in accordance with the attention weights created by the transformer layers, which are the fundamental component of Med-BERT. These attention weights are based on the results of Med-BERT. The bertviz tool⁵⁸ underwent a number of modifications and improvements so that it could more effectively display the attention patterns that were present in each layer of the pretrained model in a manner that was understandable and uncomplicated. We didn't do it for any other purpose than to enhance the general look of the site, therefore we put "SEP" tokens in the spaces between each visit. On each of the numerous layers that the model had, we saw a plethora of distinct forms of repeating patterns. Some attention heads are only present during the course of a single visit, while others point to the same codes regardless of how many visits are taken in total. These characteristics can be found in the pretrained model, which is one of the six tiers that make up the BERT transformer model. The connections between the first two layers of this model are typically syntactic. In addition to this, the pretrained model features certain attention heads that can only be activated during one visit and always direct attention to the same codes. At the middle two levels, we are beginning to see the emergence of certain patterns of attention that are medically important. These patterns take into account information that is reliant on the visit itself as well as the environment in which it occurs. Once a certain number of levels have been completed, the attention patterns become

more confused and it is hard to make sense of them. The same code may be used over several visits to depict distinct attention patterns in a variety of contexts. This serves as an example of one possible use of the code. This illustrates that Med-BERT has the capability of learning representations that are contextualized in their respective settings. The earlier code for type 2 diabetes is primarily dependent on the code for the long-term intake of insulin at the same appointment, but the later code for diabetes places more of an emphasis on the insulin code, both in the present visit and in visits from the past. The earlier code for type 1 diabetes is mostly dependent on the code for the long-term use of insulin at the same session. Because of this, one may be led to conclude that the model works out the temporal link between visits by adding segments to the timeline. Sometimes known as the "supplementary figure," has numerous additional illustrations than the previous figures. The fine-tuned model exhibits attention patterns that can be differentiated from those shown by the basic model. The fact that the fine-tuned models reflect a range of task-dependent patterns across successive layers is an indication of the generalizability and flexibility of the model for the goal of acquiring varied degrees of information in settings that are drawn from the real world. This pattern was included in the model. This pattern does an excellent job of encapsulating the strong connection that exists between diagnostic codes. Figure of the extra materials, which may be downloaded by scrolling down to this page, presents a number of other visualization patterns. These sorts of visualization patterns, in our opinion, have the potential to assist us in obtaining a better knowledge of the underlying processes that are at work inside of the neural network, which would be beneficial.

E. Pretraining cohort for the Med-BERT program

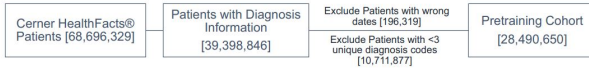
Over 600 hospitals and clinics dispersed throughout the United States make use of the de-identified electronic health record (EHR) database that is made available by Cerner Health Facts® (version 2017). It includes data on over 68 million patients, all of whom are entirely unique from one another, and it combines longitudinal data spanning the years 2000 to 2017. Information pertaining to certain patients is stored in the database. This information will be available for access at a later time. The following categories of information may be found within these records: demographics; meta-information about encounters; diagnoses; procedures; laboratory results; prescription orders; medication administration; vital signs; microbiology; surgical cases; other clinical observations; and characteristics of health systems. In order to compile the data that is shown in Health Facts®, Cerner went through the electronic medical records (EMRs) of hospitals that have entered into data use agreements with the company. Because of this, the business was able to develop the in-depth study that can now be accessed by the public. The identity of the pharmacy, clinical laboratory, and microbiological laboratory are all included in the encounter meta-information. In addition to the information on admission and payment from associated patient care facilities, this information is also included in the

encounter meta-information. In addition to that, the patient's medical history is included into the meta-information. In order to create a temporal relationship between treatment patterns and clinical information, date and time stamps are affixed to all admissions, as well as pharmaceutical prescriptions and the dispensing of those prescriptions, laboratory orders, and specimens. Because of this, it is possible to generate a picture of patient care that is more accurate.

In order to successfully deidentification of the Health Facts® database, Cerner Corporation has created operational requirements that are compliant with the Health Insurance Portability and Accountability Act (HIPAA). During the time we spent preparing the data for the pretraining, we arranged the diagnostic codes that were part of each visit in accordance with the following three criteria: The source for the diagnostic priority is the Cerner database, which has a variety of priorities for the diagnosis and is used to represent those priorities. (1) the diagnosis was recorded as being present at the time of admission; (2) the diagnosis was reported at different intervals during the visit (for example, when the patient was hospitalized); or (3) the diagnosis was only noted at the billing phase.

However, according to the findings of our investigation, including the code order throughout the whole of the assessment did not result in any appreciable leaps in performance. This, on the other hand, is not the same as randomly dispersing the codes over the test in a haphazard way.

As a consequence of this, we have arrived at the opinion that the best course of action is to publish it in this position as a stand-in, and we are going with our work on the basis that a more suitable ordering may be described at a later time. This has led us to the conclusion that the best course of action is to publish it in this area as a stand-in. persons whose medical records had less than three diagnostic codes were not permitted to take part in the study, nor were persons whose time information had been entered improperly (for instance, the discharge date was reported before the admission date). In addition, people whose medical records contained an inappropriate order of dates were also excluded from the study. When dividing the total number of patients into training, valid, and testing sets, we utilized a ratio of 7:1:2 for both the pretraining phase and the assessment phase. This ratio was also used when dividing the training sets. In the phase devoted to assessment, this ratio was also used. This ratio was calculated by taking the total number of patients from each batch and dividing it by the total number of patients. According to the information gathered, there were a total of 28,490,650 patients.



V. RESULT AND DISCUSSION:

It is particularly useful for "extreme transfer learning" paradigms, which refer to the process of fine-tuning a model with just a few hundreds of data points. The utility of Med-BERT is shown by the fact that it helps to improve prediction

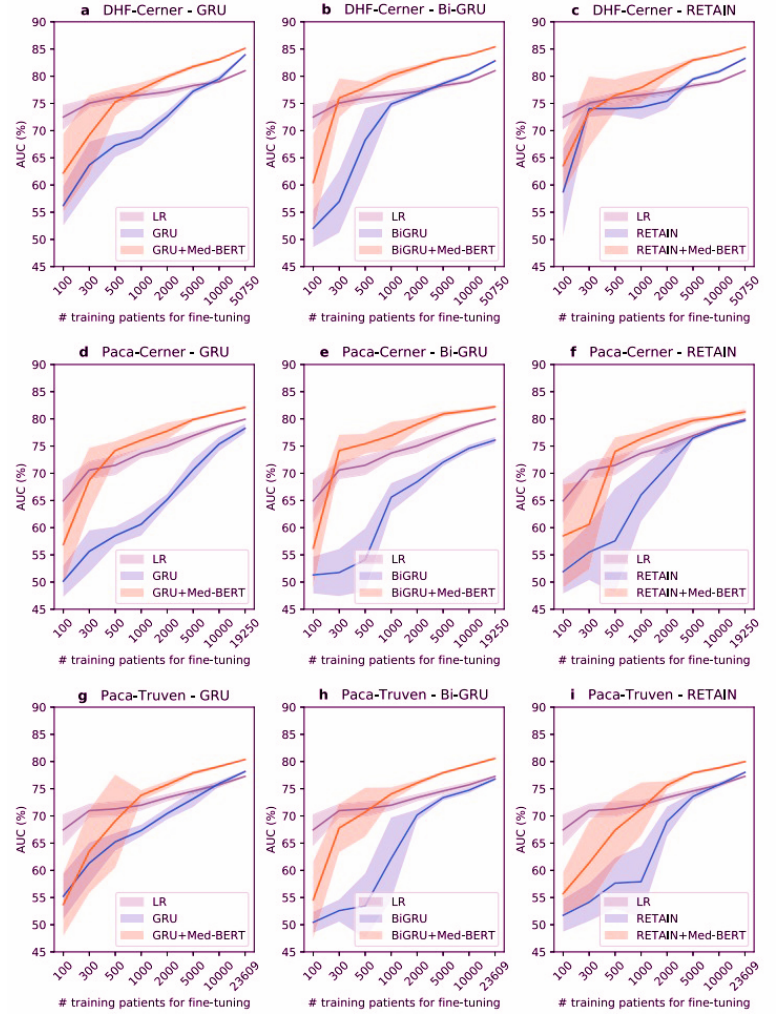


Fig. 2.

performance across a variety of jobs that each have their own unique configurations.

In order to function properly, predictive models, which often make use of deep learning, generally need at least thousands of individual data points. In order for these models to be able to manage complex cases that have never been seen before, they need to learn difficult semantics by being fed examples that indicate unique underlying sickness progressions and information about variational context. Only then will they be able to handle complex situations that have never been seen before. Nevertheless, the vast majority of deep learning algorithms are constrained in their capacity to have an in-depth knowledge of the inputs; as a result, they are unable to model the data in a manner that is comprehensive. Pretrained models have the ability to successfully overcome this challenge by utilizing more elaborate structures to more effectively capture the intricate semantics of inputs, functioning as knowledge containers, and injecting the information into new tasks. This promise can be realized via the use of more complicated structures to more effectively capture the intricate semantics

Average AUC values and standard deviations (in parentheses) for the different methods for the three evaluation tasks.			
Model	DHF-Cerner	PaCa-Cerner	PaCa-Truven
GRU	83.93 (0.13)	78.26 (0.84)	78.17 (0.21)
GRU + t-W2V	83.95 (0.24)	80.08 (1)	77.54 (0.27)
GRU + Med-BERT	85.14 (0.06)	82.13 (0.24)	80.37 (0.12)
Bi-GRU	82.82 (0.17)	76.09 (0.61)	76.79 (0.29)
Bi-GRU + t-W2V	84.23 (0.06)	79.35 (0.27)	77.44 (0.22)
Bi-GRU + Med-BERT	85.39 (0.05)	82.23 (0.29)	80.57 (0.21)
RETAIN	83.28 (0.16)	79.68 (0.32)	78.02 (0.19)
RETAIN + t-W2V	84.98 (0.02)	81.8 (0.17)	79.46 (0.18)
RETAIN + Med-BERT	85.33 (0.09)	81.3 (0.55)	79.98 (0.17)
Med-BERT_only (FFL)	85.18 (0.12)	81.67 (0.31)	79.98 (0.26)
untrained Med-BERT only	82.76 (0.13)	75.16 (0.77)	75.9 (0.18)
Logistic Regression (LR) ^a	81.01 (0)	79.94 (0)	77.28 (0)
Random Forest (RF) ^a	81.88 (0.08)	79.48 (0.31)	77.00 (0.12)

^aLR and RF input is one hot representation while other models using embeddings.
The numbers in boldface indicate the highest AUROC per task.

Fig. 3. title

of inputs. In the same way that pretrained models on other domains have been proven to be fairly effective when transferring to new tasks, our study has shown that Med-BERT is extremely helpful when transferring to new tasks by employing its bidirectional transformer, deep structure, and enormous data. our is similar to the findings that have been found for other pretrained models on other domains.

In order to simulate sequential dependencies and contextual information, Masked LM and Prolonged LOS were incorporated in the model. These two ideas were developed of and implemented into the model to be of assistance. It is possible to develop labels for one or both of these items in an unsupervised manner, which means that they may be created without any help from human annotations. The goal of Masked LM is to correctly predict a masked code by using the sequential information from both the forward and backward directions of the analysis. The purpose of the LOS analysis is to identify patients who have been linked to any visits requiring a longer-than-usual length of stay. This method uses cumulative conditions to assess if a patient has visited. If the parameter initialization in deep learning models were improved, it stands to reason that the models' performance would improve and they would arrive at a solution more rapidly. This makes intuitive sense.

On the other hand, should there be a rise in the total number of training samples, the benefits can ultimately become less substantial. Scales of samples of up to 50 and 20 K are adequate, in our opinion, when it comes to the training of excellent (converging) deep learning models, since these numbers have been determined by us. Nevertheless, after using Med-BERT, a sizable level of advancement in terms

of improvement was revealed to have occurred. For example, the RETAIN algorithm accomplishes excellent results across all three tasks; however, the inclusion of Med-BERT leads to further improvements that range from 1.62 to 2.05%. GRU and Bi-GRU's model structures are simpler than RETAIN's, hence their improvements can be greater. This raises the bar for these core models to meet or exceed RETAIN's. Also, we may conclude that Med-BERT has the potential to save researchers from the burden of constructing complex models for the numerous sickness prediction problems based on the findings of Med-BERT only, which also achieves great performance. This is feasible because of the excellent performance Med-BERT delivers.

The static embedding strategy known as t-W2V has the potential to operate as an effective performance enhancer for the basic deep learning models, in a way that is comparable to that of Med-BERT. The majority of the time, the advantages that can be obtained through the use of t-W2V are not quite as big as the advantages that can be obtained through the use of Med-BERT. One potential explanation for why t-W2V has difficulty modeling long-sequential information is that it has a shallow structure and a constrained size for the context window, neither of which can be guaranteed to perform well in all circumstances. This is because there are limits placed on the amount of space that may be used by the context window. When it is finally implemented, Med-BERT will be of significant assistance in easing the burden of data labeling in a meaningful way. This is something that may be detected by comparing the numbers of training samples required to obtain different AUC values in a given situation. Ex-3 proved that it is feasible to effectively adapt Med-BERT into a number of applications that are used for the prediction of illnesses in the real world. The bulk of the charts show that Med-BERT has the potential to greatly enhance the general performance of fundamental models when it is applied to relatively small sample numbers. For example, in the first subchart of the PaCa-Cerner chart, if we draw a horizontal line over the y-tick of 0.75, we will see that we need 1000 samples for GRU + Med-BERT, but we require nearly 10,000 samples for GRU by itself. However, if we just examine samples from the GRU, the necessary number is more than 10,000 samples. In a similar vein, we are able to demonstrate that the Bi-GRU + Med-BERT model that was trained on 5,000 samples may yield slightly superior performance than the Bi-GRU model that was trained on more than 50,000 samples alone, as shown in Supplementary Table 2A. This is possible because the Bi-GRU model was trained on more data than the Med-BERT model was. As a consequence of this, the performance of the model was improved to the point where it is now comparable to that of a training set that is over ten times larger. Because we will be utilizing Med-BERT, we will be able to gain considerable cost savings in respect to the cost of data collection, which may at times be very dear. This cost reduction will be obtained in relation to the more than 9000 samples we will be analyzing. Even before collecting a significant number of annotated samples, researchers and medical professionals

in this circumstance may use Med-BERT to quickly get a knowledge that is both extensive and acceptable about the progressions of unique disorders. Even though Med-BERT allows deep learning models across the board for all of the training sample sizes that were assessed, it is nevertheless the case that models powered by Med-BERT do not exceed the non-deep learning baseline model LR for the lowest training sample sizes (n fewer than 500). This is true even though Med-BERT enables deep learning models across the board for all of the training sample sizes that were evaluated. This is in agreement with what has been mentioned in the research literature, which shows that LR continues to be a competitive predictive model despite having reduced training sample sizes in a number of studies¹⁴. This is because this is in line with what has been stated in the research literature. The benefit of LR is that it has a simple and shallow structure, which makes it much easier to fit the data based on even a small number of samples. This is because the structure is so simple and shallow. In contrast to this are the deep learning models, which have a convoluted structure and a huge parameter space to work with. On the other hand, when the number of people participating in the training increases, the magnitude of this advantage will diminish to some extent. Because of this, in practice, we would recommend using Med-BERT fine-tuning for the situations in which the training sample size is sufficiently large (for example, $n \geq 500$). In the most current release of Med-BERT, the ICD-9 and ICD-10 codes have been merged into a single vocabulary, bringing the total number of tokens in that release to 82,000. Our vocabulary, in contrast to BEHRT and G-BERT, covers a bigger variety of subject areas and may be used in a higher number of settings. Additionally, our vocabulary is more versatile. We are of the opinion that the International Classification of Diseases (ICD), which is a worldwide health information standard that has been endorsed by the World Health Organization and is utilized by numerous organizations from over one hundred countries all over the world, will significantly assist in the model's ability to be transferred to other locations. This was shown by our PaCa-Truven research, in which we investigated the efficacy of our models by making use of a cohort that was created from a dataset that included information on health insurance. The data for this cohort came from a database. For the sake of this investigation, the EHR modality that we decided to use was BERT, which is a sophisticated contextualized embedding technique in NLP. This was done so that we could get the most accurate results. However, there are various methods of thinking, such as ULMFiT⁴⁶, ELMo²⁶, GPTs^{27,28,59}, and so on. Some of them are included here. It is very likely going to be necessary to do an examination of all of these different pretraining and fine-tuning alternatives for electronic health records (EHR). We are going to save it for use in tasks that may come up in the future. The work that is now provided suffers from a number of deficiencies that absolutely have to be fixed. In the first step of this process, we limited ourselves to just utilizing the information on the diagnosis that was supplied in ICD format. Second, we did not take into

consideration the length of time that had gone in between each participant's visits in this study. As a consequence, there is a possibility that some temporal information was lost as a result of this oversight. Third, we did not conduct an exhaustive investigation into the sequence of ideas included in each visit, and the current setup, which is based on code priority, may not be stable enough for long-term use. In the future, further research will be done on the design of a number of different pre-training activities, and there will also be testing done on a variety of other kinds of fine-tuning tasks that go beyond sickness prediction. Both of these things will take place in the future. In addition, we wish to include more inputs into the Med-BERT model, such as time, medications, procedures, and laboratory tests. These are all examples of potential inputs. In the near future, we also intend to study task-specific visualizations and interpretations.

In conclusion, we presented the contextualized embedding model known as Med-BERT. This model had been pretrained on a considerable quantity of structured EHR data, and we went on to test it in sickness prediction tasks. We developed assignments and input formats that were specifically tailored to a certain topic and were taught earlier. Extensive analysis has shown that Med-BERT may help improve the prediction performance of baseline deep learning models on a wide range of training sample sizes. These findings were based on the findings of the research described above. This and similar experiments demonstrated Med-BERT's potential to enhance the prediction performance of standard deep learning models. Because of the visualization module, we were able to dive more thoroughly into the semantics of the data as well as the working processes of the model. This gave us possibilities to recognize notable examples in both of these areas. After that, clinical professionals validated these samples, which demonstrates that Med-BERT is capable of capturing the semantics across EHRs during both the pretraining and finetuning phases of the process. From a methodological point of view, our research demonstrates how contextualized embedding of structured EHR data may be carried out, therefore validating its feasibility and demonstrating its potential utility. In practice, our pretrained model makes it feasible to construct successful deep learning prediction models using just a small number of data points from the training set. This is made possible by the fact that it only requires a limited number of data points to be provided by the training set.

VI. EVALUATION OF PROPOSED MODEL

These studies were carried out after we analyzed two separate illness prediction tasks using three different cohorts derived from two different databases. There are two projects that need to be completed, and they are DHF and PaCa. Both of these tasks were completed with the assistance of Cerner, which led to the establishment of the DHF-Cerner and PaCa-Cerner cohorts. Because we only wanted to use Truven for the pancreatic cancer prediction test, we were able to create the PaCa-Truven cohort. This allowed us to evaluate the extent to which our findings were applicable to a wider

population. The comprehensive definitions of the cohorts are laid out in detail in the "Methods" section of the report. Our definition of illness prediction tasks is more complicated because it needs the phenotyping from several viewpoints. These views include the presence of particular diagnostic codes, medicinal prescriptions, procedures, laboratory test results, and sometimes even the frequency of occurrences within preset time intervals. This is in contrast to BEHRT and G-BERT, whose assessment tasks merely consist of the prediction of particular codes and are hence analogous to the tasks that are performed during pre-training. Because of this, we claim that the evaluation tasks that we utilize are more advantageous to the process of assessing the generalizability of Med-BERT and are more appropriate to the real world (in contrast to BEHRT). Experiments were carried out in accordance with one of the following categories for each of the three tasks: (1) to figure out how much Med-BERT can add to the current best methods. (2) Evaluate the similarities and differences between Med-BERT and t-W2V, a state-of-the-art static clinical word2vec-style embedding that was also trained on the whole Cerner cohort. (3) to determine how much the pretrained model helps transfer learning with varied training sample sizes. For each fine-tuning job, a random subset of the original group was chosen, and that subset was then split into training, confirmation, and testing sets in the ratio 7:1:2. Since we had a lot of patients who weren't in the training set, we gave priority to putting samples in the test set. This was to make sure that our test sets didn't have any patients who had been in the Med-BERT training set before. This was done in order to ensure that our test sets accurately reflected the population of patients who were being tested. We decided to use the area under the curve (AUC) as our primary measure of evaluation so that we could evaluate performance more accurately. This assessment measure has had significant application in a range of past studies dealing with the forecasting of illnesses. Other information on other performance evaluation metrics may be found in the tables 1 and 2 of the supplemental materials. To measure how much more effective modern models become when layered with pretrained Med-BERT, we simply compare the performances of the base models with the performances of the base models layered with Med-BERT. This allows us to determine whether or not pretrained Med-BERT has boosted the power of state-of-the-art base models. This makes it possible for us to establish whether or not the incorporation of Med-BERT into the basic models results in an increase in the power of the models. The GRU53, Bi-GRU54, and RETAIN12 architectures are the ones that we use most often for our core recurrent neural network (RNN) models. We also added RETAIN, which is a well-known sickness prediction model that makes use of double GRUs in combination with attention. This was done despite the fact that it had been shown that GRUs are incredibly competitive baseline models. In addition, we demonstrated the results by using just Med-BERT; more specifically, we included only FFL in the very last layer of Med-BERT. This Med-BERT-only model will offer an evaluation that exceeds what is possible with RNN-based

models due to the model's exclusive use of the Med-BERT. We also compare the performance of the pretrained Med-BERT architecture to the performance of the untrained Med-BERT architecture in order to examine the influence that pretraining has while dealing with vast volumes of data. As a matter of completeness, we also included L2-regularized logistic regression (L2LR) and random forest (RF) as baseline models. L2LR and RF are well-known methods that do not incorporate deep learning. Both of these approaches make use of the conventional multishot input format as their source of data. In order to analyze the differences and similarities between Med-BERT and static embeddings, we choose to conduct Ex-2 using the t-W2V model as our research tool. The results of a prior study⁴⁵ served as the basis for our decision to make use of t-W2V as a representation for non-contextualized static embeddings. In this work, various different static embedding approaches, such as word2vec²⁴, fasttext⁵⁵, and pointwise positive mutual information singular value decomposition⁵⁶, were compared with one another. According to the findings of this comparison, t-W2V performed the best out of all the candidates in the analyzed illness prediction job. It should be pointed out that Glove²⁵ is a competent alternative to word2vec (w2c) for the embedding of static EHR ideas, and it was reported as having a performance that was comparable to that of w2c. This is something that should be taken into consideration. Because of this, and in the interest of keeping things as simple as possible, we have settled on t-W2V as the baseline for our static embedding. In the case of Example 3, in order to determine the value-added of Using Med-BERT with a variety of fine-tuning training sizes, we selected training data samples for fine-tuning with increasing sizes for each cohort. We did this so that we could evaluate the outcomes side by side. Intuitively, it looks as if the pretrained model would be of more aid when the training size is smaller since it helps infuse a larger variety of information. This is because it allows for more accurate predictions to be made. We supplied the average area under the curve (AUC) as well as the standard deviation for each model. These values were based on ten runs with randomly initialized prediction head weights, and they were calculated for Ex-1 and Ex-2, where the whole fine-tuning training cohort was used. These findings are an inference made from the information gathered during the training phase of the project. Following the calculation of the average area under the curve (AUC) as well as the standard deviation for each cohort, we carried out a random bootstrap sampling ten times over all of the iterations in Ex-3.

VII. CONCLUSION

To conclude, it can be said that in the rapidly evolving field of natural language processing (NLP) in healthcare, standardized and trustworthy datasets are essential for the development and evaluation of algorithms and models. The Real-MedNLP dataset was created to address the lack of anonymized medical text data in languages other than English, and it contains real clinical data from case and radiology reports in both Bangla and English translations. The dataset's emphasis on

practical tasks and clinical value makes it a valuable resource for future medical NLP research aimed at developing models suitable for use in actual healthcare environments. Overall, the construction of the Real-MedNLP dataset fills a significant need in the study of medical NLP and paves the way for the advancement of NLP methods in healthcare. In future, we will try to include Bengali medical text documents in our dataset to make our research more effective and fruitful.

REFERENCES

- [1] Kim, Y., Kim, J., Lee, J. M., Jang, M. J., Yum, Y., Kim, S., Shin, U., Kim, Y., Joo, H. J., Song, S. (2022). A pre-trained BERT for Korean medical natural language processing. *Scientific Reports*, 12(1). <https://doi.org/10.1038/s41598-022-17806-8>
- [2] MMBERT: Multimodal BERT Pretraining for Improved Medical VQA. (2021, April 13). IEEE Conference Publication — IEEE Xplore. <https://ieeexplore.ieee.org/document/9434063>
- [3] Juric, D., Stoilos, G., Melo, A., Moore, J., and Khodadadi, M. (n.d.). A system for medical information extraction and verification from unstructured text. *Proceedings of the AAAI Conference on Artificial Intelligence*. Retrieved April 8, 2023, from <https://doi.org/10.1609/aaai.v34i08.7042>
- [4] Pandey, B., Pandey, D. K., Mishra, B. K., Rhmann, W. (2021). A comprehensive survey of deep learning in the field of medical imaging and medical natural language processing: Challenges and research directions. *Journal of King Saud University - Computer and Information Sciences*, 34(8), 5083–5099. <https://doi.org/10.1016/j.jksuci.2021.01.007>
- [5] NLP-Based Prediction of Medical Specialties at Hospital Admission Using Triage Notes. (2021, August 1). IEEE Conference Publication — IEEE Xplore. <https://ieeexplore.ieee.org/abstract/document/9565791>
- [6] Soufyane, A., Abdelhakim, B. A., Ahmed, M. H. (2021). An Intelligent Chatbot Using NLP and TF-IDF Algorithm for Text Understanding Applied to the Medical Field. *Advances in Science, Technology Innovation*, 3–10. https://doi.org/10.1007/978-3-030-53440-0_1
- [7] Liu, L., Grossman, R. H., Mitchell, E. G., Weng, C., Natarajan, K., Hripcsak, G., Vawdrey, D. K. (2021). A deep database of medical abbreviations and acronyms for natural language processing. *Scientific Data*, 8(1). <https://doi.org/10.1038/s41597-021-00929-4>
- [8] Roy, A., Pan, S. (2021). Incorporating medical knowledge in BERT for clinical relation extraction. *Empirical Methods in Natural Language Processing*. <https://doi.org/10.18653/v1/2021.emnlp-main.435>
- [9] August, T. (2022b, February 28). Paper Plain: Making Medical Research Papers Approachable to Healthcare Consumers with Natural Language Processing. *arXiv.org*. <https://arxiv.org/abs/2203.00130>
- [10] Zeng, J., Rubin, D. L., Henry, S., Wood, D. E., Shachter, R. D., Gensheimer, M. F., Rubin, D. L. (2021). Natural Language Processing to Identify Cancer Treatments With Electronic Medical Records. *JCO Clinical Cancer Informatics*, 5, 379–393. <https://doi.org/10.1200/cci.20.00173>
- [11] Jiang, F. et al. Artificial intelligence in healthcare: past, present and future. *Stroke Vasc. Neurol.* 2, 230–243 (2017).
- [12] Yu, K.-H., Beam, A. L., Kohane, I. S. Artificial intelligence in healthcare. *Nat. Biomed. Eng.* 2, 719–731 (2018).
- [13] Chen, M., Hao, Y., Hwang, K., Wang, L. Wang, L. Disease prediction by machine learning over big data from healthcare communities. *IEEE Access* 5, 8869–8879 (2017)
- [14] Wang, H. et al. Predicting hospital readmission via cost-sensitive deep learning. *IEEE/ACM Trans. Comput. Biol. Bioinforma.* 15, 1968–1978 (2018).
- [15] Davenport, T. Kalakota, R. The potential for artificial intelligence in healthcare. *Future Healthc. J.* 6, 94 (2019).
- [16] Lysaght, T., Lim, H. Y., Xafis, V. Ngiam, K. Y. AI-assisted decision-making in healthcare. *Asian Bioeth. Rev.* 11, 299–314 (2019)
- [17] Ahmed, Z., Mohamed, K., Zeeshan, S. Dong, X. Artificial intelligence with multifunctional machine learning platform development for better healthcare and precision medicine. *Database* 2020, baaa010 (2020). <https://doi.org/10.1093/database/baaa010>
- [18] Manogaran, G. Lopez, D. Health data analytics using scalable logistic regression with stochastic gradient descent. *Int. J. Adv. Intell. Paradig.* 10, 118–132 (2018).
- [19] Keerthika, T. Premalatha, K. An effective feature selection for heart disease prediction with aid of hybrid kernel SVM. *Int. J. Bus. Intell. Data Min.* 15, 306–326 (2019).
- [20] Sadek, R. M. et al. Parkinson's disease prediction using artificial neural network. *Int. J. Academic Health Med. Res.* 3, 1–8 (2019).
- [21] Payan, A. Montana, G. Predicting Alzheimer's disease: a neuroimaging study with 3D convolutional neural networks. Preprint at <http://arxiv.org/abs/1502.02506> (2015).
- [22] Choi, E. et al. RETAIN: An Interpretable Predictive Model for Healthcare using Reverse Time Attention Mechanism. *Adv. Neural Inf. Process. Syst.* 29, 3504–3512 (2016)
- [23] Choi, E., Bahadori, M. T., Schuetz, A., Stewart, W. F. Sun, J. Doctor AI: Predicting Clinical Events via Recurrent Neural Networks. In *Machine Learning for Healthcare Conference*, 301–318 (MLHC, 2016).
- [24] Rajkomar, A. et al. Scalable and accurate deep learning with electronic health records. *NPJ Digital Med.* 1, 18 (2018).
- [25] Esteva, A. et al. Dermatologist-level classification of skin cancer with deep neural networks. *Nature* 542, 115–118 (2017).
- [26] Poplin, R. et al. Prediction of cardiovascular risk factors from retinal fundus photographs via deep learning. *Nat. Biomed. Eng.* 2, 158 (2018).
- [27] Coudray, N. et al. Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning. *Nat. Med.* 24, 1559–1567 (2018).
- [28] Chung, S. W. et al. Automated detection and classification of the proximal humerus fracture by using deep learning algorithm. *Acta Orthop.* 89, 468–473 (2018).
- [29] Shrivastava, A., Singh, S. Gupta, A. In *Proceedings of the IEEE International Conference on Computer Vision*, 843–852.
- [30] Cho, J., Lee, K., Shin, E., Choy, G. Do, S. How much data is needed to train a medical image deep learning system to achieve necessary high accuracy? Preprint at <https://arxiv.org/abs/1511.06348> (2015).