

Enhancing Predictive Accuracy and Resource Utilization in Disease Prediction: A Comparative Study of Small Versus Large Language Models

V Karthick, Associate Professor
Department of CSE
Rajalakshmi Engineering College
Chennai, India
vkarthick86@gmail.com

Tamanna, UG Student
Department of CSE
Rajalakshmi Engineering College
Chennai, India
210701281@rajalakshmi.edu.in

Sudhir S, UG Student
Department of CSE
Rajalakshmi Engineering College
Chennai, India
210701269@rajalakshmi.edu.in

ABSTRACT - This project involves a comparative analysis of small language and large language models where it evaluates the models' capabilities in deciphering complex relationships between diseases and their symptoms using the distilgpt2(SLM) and gpt2(LLM) models. We utilized a large dataset where both small and large language models were used to study diseases and their symptoms. The models were evaluated based on their computational efficiency, predictive performance and the resources required for the operation. The methodology includes loading datasets, tokenization and model setup, monitoring training and validation losses, hyperparameter tuning to optimize the models and eventually generating text. The study showed SLM's proficiency in producing context-aware responses which showed an inference time of 0.008518 secs whereas LLM shows its strength in generating refined, comprehensive text with a larger inference time of 0.015379 secs. The outcome showed the ability of large language models to possess higher predictive accuracy and they demanded significantly higher computational resources as compared to small language models, which may not be advised to use in resource-constrained environment. In contrast, small language models have efficient resource usage despite their low accuracy make their application more pertinent where computational resources are limited. This study paves the way for understanding the blueprint for health informatics practitioner highlighting the importance of domain-specific training in enhancing predictive accuracy and resource

utilization of language models. It also underscores the requisite for a new approach to deployment of language models in healthcare settings.

Keywords - *Small Language Models, Large Language Models, BERT, Tokenizer, DistilBERT, resource-constrained, disease prediction*

I. INTRODUCTION

The area of investigation of small and large language models in disease prediction is a rapidly growing field of study that has significant potential for the future of healthcare. The SLMs with their few parameters and smaller datasets are characterized by their speed and efficiency in processing, thus, they are suitable for the applications where computational resources are limited making them suitable for real-time applications. The other way around, large language models (LLMs) are characterized by their huge number of parameters and the intensive training on the large datasets, thus, they can perform the complex language tasks with high accuracy and are capable of understanding complex symptom descriptions and providing more accurate disease associations. The focus in this project is on utilizing a dataset from Hugging Face to predict disease based on reported symptoms and comparing the outcomes from small and large language models. This involves leveraging these language models to understand and interpret symptoms data and accordingly map them to potential diseases.

The methodology involved in the project includes feeding the language models with symptom data and analyzing the outputs to determine the most likely disease outcomes. The goal is to make most out of the SLMs and LLMs to create a resilient system that can help healthcare providers in the initial diagnosis of diseases, thus, the diagnostic process is going to be made easier. This study helps us gain deep insights on practicality of introducing such language models in the domain of healthcare considering factors like accuracy, computational speed and efficiency of prediction. The outcomes of the research will be instrumental in determining the most effective way to implement language models for medical help, which will, in turn, help in the betterment of patient care and resource management in the healthcare sector.

Approach: Understanding SLMs and LLMs

Small Language Models (SLMs):

- **Architecture:** Language Models with a reduced context are called Small Language Models (SLMs). SLMs describe a neural network with fewer layers. They often utilize less numbers of layers and parameters, thus ensuring faster training and inference time.
- **Training:** SLM that are built by only using smaller datasets are less computationally expensive. This training method is optimal for situations, where there is no data available in high volume to work with.
- **Efficiency:** The lean construction or architecture of the SLMs is what leads to fast operational times and maximum efficiency[9]. They can process language tasks rapidly, making them ideal for real-time applications.
- **Application:** Due to their resource efficiency, SLMs are being effectively implemented in situations where computational resources are a limitation or a real-time analysis is necessary.

Large Language Models (LLMs):

- **Architecture:** Large Language Models (LLMs) are based on a neural network of a great complexity that involves a deep learning approach. They are characterized by a higher

number of trainable parameters that give them the ability to express complex patterns in data.

- **Training:** LLMs are trained with the aid of massive data that has an enormous variety. This procedure is an essential but energy-prodigious part of the model that enables it to carry out complicated language tasks.
- **Predictive Performance:** LLMs are well-known for their high accuracy and generating output that is not only detailed but also of good quality. Such language mastery helps to establish disease correlations with greater precision.
- **Resource Requirements:** The elaborate features of LLMs bear the tradeoff of more computing resources. They need more powerful hardware and consume more energy, thus in the places of resource scarcity, their performance can be lower.

II. LITERATURE REVIEW

[1] Victor Sanh, Lysandre Debut, Julien Chaumond, Thomas Wolf in this paper attempted to create a small language model that is equally efficient without compromising its performance. It mentions the issues they usually face while training large language models when the resources and data available are limited. Their approach of using knowledge distillation during the pre-training phase is notable and shows promising results. Their work provides a motivation to develop more efficient small language models with high accuracy and better performance.

[2] In this paper, the authors are proving the efficiency of BERT in different NLP cases, such as question answering and language inference tasks. The BERT model was put to test and it was observed how it outweighed the performance of other models. BERT's capability to perform complex language processing tasks with utmost efficiency eventually leads to the understanding that it can similarly perform disease prediction with high accuracy by interpreting the intricate details of connection between symptoms and diseases.

This study[3] by Timo Schick and Hinrich Schütze analyzed the performance of small language models

compared to large language models. They showed how small models trained with few parameters as compared to large models can perform equally in predicting the result. By converting textual inputs into cloze questions and leveraging gradient-based optimization, these “greener” models achieve impressive natural language understanding

This research by Leonardo Ranaldi and André Freitas in [4] tried to develop a method to bridge the gap in reasoning skills between small and large language models. With the help of Instruction-tuning-CoT method they provide SLMs with ability to perform multi step controlled reasoning when elicited with the CoT mechanism.

This paper [5] by Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, and others show how scaling up the models improve their performance sometimes even rivaling prior state-of-the-art fine-tuning approaches. They trained GPT-3 with 175 billion parameters and tested their performance without any fine tuning. They demonstrated how GPT-3 produced strong results on translation, question-answering and other reasoning based challenges.

[7] The paper "Large language models in health care: vulnerabilities, risks, and challenges" discusses the LLMs opportunities in healthcare. It briefly touches upon the problems of using machine learning systems in biomedical and clinical fields being applied in clinical environments, and those challenges which need to be solved in order to widely adopt their usage in this field.

This paper introduces us to TinyLlama which is a compact language model pre-trained on approximately 1 trillion tokens for around 3 epochs [8]. It significantly outperforms existing open-source language models with comparable dimensions.

In [9] Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, Xinrong Zhang, Zheng Leng Thai, Kaihuo Zhang, Chongyi Wang, Yuan Yao, Chenyang Zhao, Jie Zhou, Jie Cai, Zhongwu Zhai, Ning Ding, Chao Jia, Guoyang Zeng, Dahai Li,

Zhiyuan Liu, and Maosong Sun explores the potential of Small Language Models as resource-efficient alternatives to Large Language Models. While focusing on SLMs, our approach exhibits scalability in both model and data dimensions for future LLM research.

In [10], Awobade, Oduwole, and Kolawole investigate the effect of compression techniques on small language models. Practical strategies such as pruning, knowledge distillation, and quantization have been applied on AfriBERTa, a low-resource and small-data BERT system for Koine Greek. These techniques were performed to find out the which one had more effect in improving the model's efficiency

This paper [11] showed NLPs effectiveness varies across diseases. They conducted a search across major databases including PubMed, Embase, Web of Science, and Scopus, up to December 2023, using keywords related to NLP, LLM, and infectious diseases.

In [12] the research explores the use of LLMs and deep learning techniques for predicting disease from symptoms. It analyzed two Medical Concept Normalization—Bidirectional Encoder Representations from Transformers (MCN-BERT) models and a Bidirectional Long Short-Term Memory (BiLSTM) model and demonstrated how accurately they predict the output.

The research in [14] by Minki Kang, Seanie Lee, Jinheon Baek, Kenji Kawaguchi, and Sung Ju Hwang devised a new approach to performing knowledge intensive tasks in small language models. By leveraging the use of external-knowledge resources, they utilize knowledge-augmented reasoning distillation to improve the performance of SLMs.

This paper [16] focuses on developing open-source language models in the medical application area. It includes text collection from PMC articles, language model's calibration, and implementing specific optimization strategies.

This paper by Wang et al. introduced a novel framework called Step-by-step knowLedge dIstillation

fraMework for recommendation (SLIM)[17]. In this research paper they tried to elevate the reasoning capabilities of large language models to work in a resource efficient manner by enabling sequential recommenders. They distill knowledge from large models and feed it into small models based on user behavior sequences.

In this paper [18] addresses the concerns associated with developing large language models with up to trillion parameters. To address these limitations, they present MiniCPM which consists of: 1. 2B and 2. 4B non-embedding parameter variants. It is worth to note that in each of these categories these Small Language Models (SLMs) themselves perform very well and are on a par with 7B-13B LLMs.

In [19] they introduced LLaVa-Phi which is a multi-modal assistant that harnessed the power of small language models Phi-2. It showed that even small language models even those with 2.7B parameters can excellently engage in complex dialogues that combine text and visual elements trained on quality datasets.

This research paper [20] analyzes how well pruning, knowledge distillation, and quantization work on small language models that are low on resources. This study provides that compression techniques eventually improve SLMs efficiency.

The research gap in this study is the absence of the comparative analysis of small and large language models (SLMs and LLMs) majorly in disease prediction in terms of their predictive accuracy and resource utilization. Even though there are studies on the individual capabilities of SLMs and LLMs, this research aims to bridge the gap by giving a direct comparison within the stringent constraints and precise requirements of healthcare informatics.

The aim of the study is to:

1. To find out the linkage between diseases and symptoms by finding out how efficient their response is and its predictive performance.
2. The resources needed by language models to carry out the functions needs to be specified.
3. In order to improve the predictive performance of

the model, it is necessary to understand the magnitude of domain-specific training.

4. To understand the practicality of these models in real world environments, it is mandatory to deploy and put them to use.

III. MATERIALS AND METHODS

- Datasets - A record of symptoms and diseases associated with those, possibly from medical records or public health databases.
- Preprocessing Tools - the data used by the models needs to be cleaned and formatted by software for analysis, such as tokenization and normalization.
- Language Models - DistilGPT-2 like small language models and GPT-2 like large language models

SOFTWARE REQUIREMENTS

This project requires a system running Windows 11 and Python version 3.8 or later. You'll also need several Python libraries for various tasks: transformers and torch for handling deep learning models, numpy for numerical computations, matplotlib and pandas for data visualization and manipulation, and scikit-learn for additional machine learning tools. Notably, pandas will be used for data cleaning and preparation, while scikit-learn can provide supplementary machine learning algorithms beyond the deep learning models.

HARDWARE REQUIREMENTS

The hardware requirements for SLMs include Standard CPUs or entry-level GPUs would suffice and for LLMs are High-performance GPUs with substantial VRAM, or cloud-based solutions like Google Colab Pro or AWS for extensive training.

IV. EXISTING SYSTEM

The existing system involves using GPT-2 which is a large language model for prediction of disease based on the symptoms and

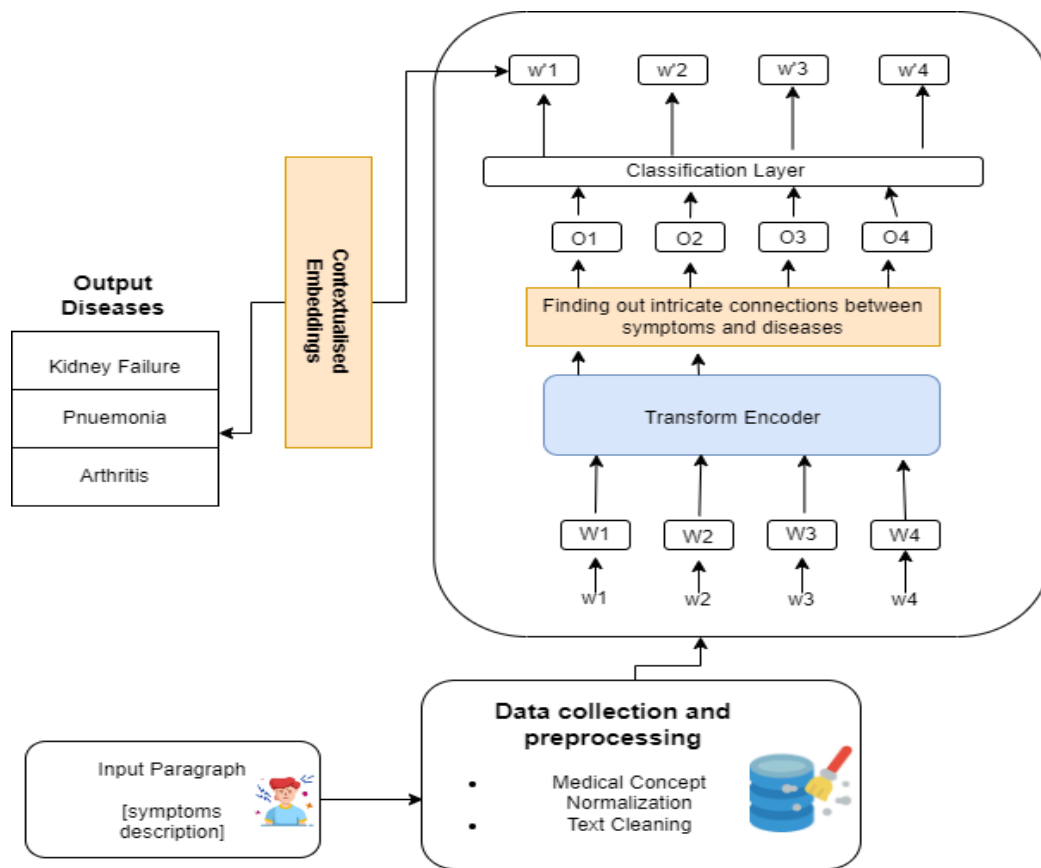


Fig.1. Architecture Diagram

vice versa. Through bespoke fine-tuning they have been trained on domain specific data to find the diseases based on the input symptoms [7]. They were trained on massive datasets to capture patterns and apply the medical knowledge. A well-known algorithm in this area is the Medical Concept Normalization-Bidirectional Encoder Representations from Transformers (MCN-BERT), which is a very accurate tool for disease prediction from symptom descriptions [12]. They tackled this challenge by harnessing the power of language models like BERT that can learn by finding out intricate connections between the symptoms and the diseases. This model can enhance their medical data thereby improving domain-specific knowledge. Bidirectional LSTM layers were studied to investigate the influence of their use on the comprehension of the contextual relations

between symptoms and diseases. The technique of hyperparameter optimization using Hyperopt was employed to increase the model's performance and thus, the model will be able to generalize well to the new data. In the existing system MCN-BERT model was used that followed the below steps:

- **Data collection and preprocessing** - This step involved collecting input datasets which consisted of symptoms mapped to their corresponding disease labels. The data was preprocessed by using a medical tokenizer where symptoms were diligently tokenized to enhance the contextual understanding of medical terms.
- **BERT model and tokenizer initialization** - The model was fed with pre-trained model weights which contain detailed linguistic knowledge that enhance the model's ability to

grasp intricate medical terms. BERT models were architected in such a way to include multiple transformer encoder layers and a specialized tokenizer was initialized. This tokenizer seamlessly converted the symptom description to tokens that were in turn integrated into the model as input.

- **Model Training** - The course of MCN-BERT training includes a thorough analysis of all stages to guarantee the correct learning process. Firstly, batches of tokenized symptom description and corresponding disease labels are organized as in the equation below

$BatchData = PrepareBatches(TokenizedSymptoms, DiseaseLabels)$

where BatchData and PrepareBatch are the prepared training batches, and function for PrepareBatches orchestrates this preparation. The following is the Model Forward Pass in which we propagate the batches through the BERT model using its forward pass mechanism.

- **Model Evaluation** - Well-known evaluation metrics, including accuracy, precision, recall, and F1-score, were used to assess the models' classification performance.

$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$

$Precision = \frac{TP}{TP + FP}$

$Recall = \frac{TP}{TP + FN}$

$F1 - score = 2 * (10) \left(\frac{Precision * Recall}{Precision + Recall} \right)$

V. PROPOSED SYSTEM

The proposed work aims to perform a comparative research on small language models (SLMs) and large language models (LLMs) in a setting of the disease prediction. This work is aimed at making an achievement in the existing knowledge[14] about disease prediction with a comparative examination of small language models (SLMs) and large language models (LLMs). The aim is to assess and compare the forecasting ability, workload, and resource deployment of both types of models.

The research will further utilize these SLMs such as DistilGPT-2 and LLMs which includes the GPT-2 or BERT which have been shown to perform impressively in natural language processing tasks in diverse ways. These models will be adjusted to the data set of symptoms and illnesses, their performance will be evaluated based on the level that they can correctly predict diseases from the symptom descriptions. Similarly, the study will evaluate both models for inference time in addition to the predictive accuracy parameter. This approach will help to determine a comparative study of resource usage by both SLMs and LLM. Besides, the study will use the comparison of the loss function of the both models to measure their predicting accuracy level.

The proposed work will, therefore, prove to be a vital asset in the field of disease prediction as it will offer an exhaustive appraisal of SLMs and LLMs. The findings from this experiment might serve as a guide for later studies in the field and help create better and comprehensive disease prediction models. This work is an innovative and well separated concept that not only considers the precision of the models but also the computational efficiency and the resource utilization. This makes it a very important contribution to the field of disease prediction.

VI. METHODOLOGY

The methodology of this study is designed to provide a comprehensive comparison of Small Language Models (SLMs) and Large Language Models (LLMs) in the context of disease prediction. It aims to evaluate and contrast the predictive performance, computational efficiency, and resource utilization of both types of models. The steps involved include:

- **Data Preprocessing:** Like the present one, the data is prepared for processing by standardizing the symptoms in the form of a single string separated by commas.
- **Dataset Preparation:** The processed data is then

converted in a form that can be used for model training. This is tokenization of data input and the creation of the data loaders for training.

- **Model Initialization:** The first thing that is done with evolved pre-trained LMs and evolving SLMs is the initialization. For SLMs, for instance, models like DistilGPT-2 can be toyed with, and for LLMs, GPT-2 or BERT can be used. The models are initialized with pre-weighted parameters and relocating these weights to the required device (CPU or GPU).
- **Model Training:** The layers are trained by using a regular training loop. During every epoch, the parameters of the model are adapted to minimize the loss value. Perturbation delivered through the forward method of the model provides the loss value to the output in the loss attribute.
- **Model Evaluation:** After training, model accuracy is tested on the validation set considered. The way of the performance of the network is tested as the loss function is used which shows the difference between the predicted values by model and actual data. The metrics achieved in the present model (precision, accuracy, recall, and the F1-score) can be calculated in order to conduct a thorough assessment.
- **Inference Time Measurement:** Both such models' inference time is measured. It meant, feeding an input text into both the models and measuring the time, how fast these models generated the output.
- **Inference Time Comparison:** Two models are compared in terms of the inference times and a corresponding bar chart is used to illustrate this.
- **Validation Loss Comparison:** Bar graphs are employed to show the validation losses of the two models under testing. Moreover, the losses are compared.

VII. RESULT

The results highlight the trade-offs between SLMs and LLMs. While LLMs may provide better performance on the training data, SLMs can often generalize better to unseen data and provide faster inference times. This makes SLMs a viable option for applications where computational resources are limited or a fast response is required.

In this comparison:

Training Loss: The SLM has a higher training loss difference than the LLM. It could be that the SLMs simplifying of the architecture and the few parameters involved does not enable as close the fitting of the training data as can be done by the LLM.

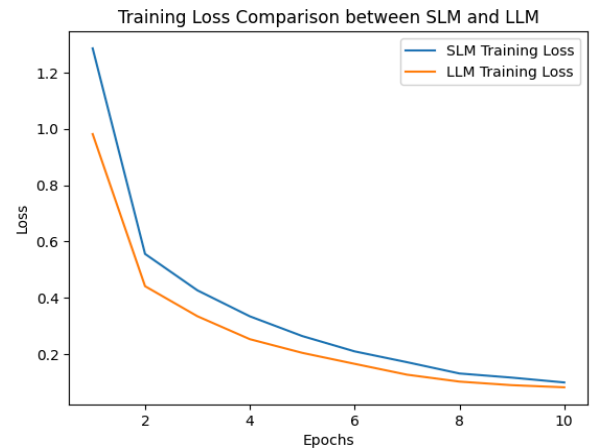


Fig.2. Training Loss Comparison

Validation Loss: In comparison, the validation loss of the SLM is less than that of the LLM. This demonstrates that the SLM has good generalizability to the data that the neural network has not encountered despite its higher training loss. This could be in connection with the discriminatory power of SLMs, where by avoiding data overfitting, they evade the overfitting problem of larger models such as the LLMs.

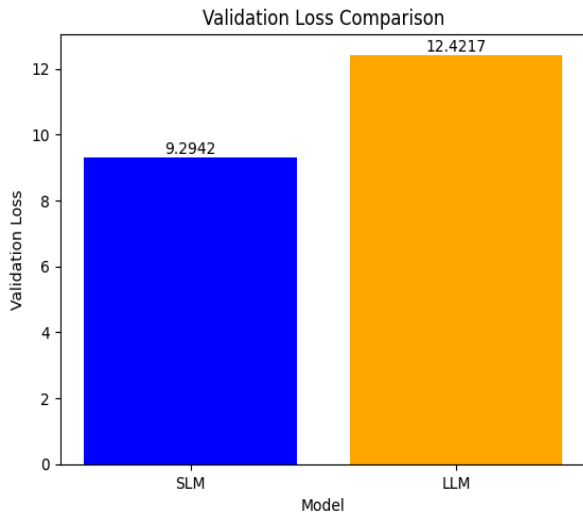


Fig.3. Validation Loss Comparison

Inference Time: The SLM speeds up the inference process relative to the LLM. This was to be anticipated as the SLMs are usually more computationally efficient owing to their small size/dimensionality and lesser number of parameters.

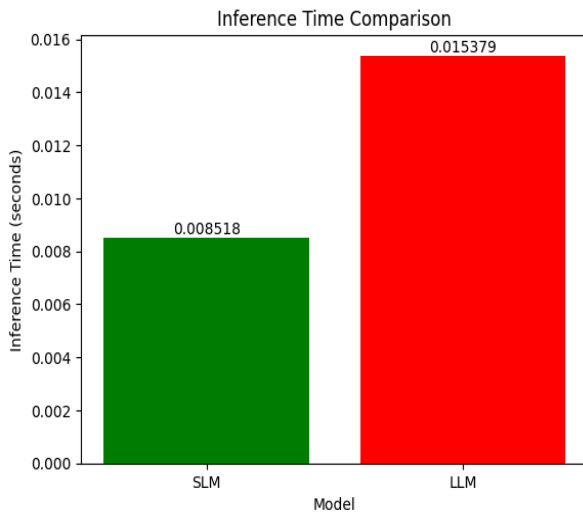


Fig.4. Inference Time Comparison

| Model | Training Loss | Validation Loss | Inference Time |
|-------|---------------|-----------------|----------------|
| SLM | Higher | 9.2942 | 0.008518 secs |
| LLM | Lower | 12.4217 | 0.015379 secs |

Table1. Model Comparison

VIII. DISCUSSION

Disease prediction is a critical field that requires comprehensive comparisons between Small Language Models (SLMs) and Large Language Models (LLMs), as this study provides. The research shows the interplay between SLMs and LLMs. Actually, despite possible superiority of LLMs on training data, the corresponding feature of SLMs generalizing better to the unseen data and offering faster inference time remains[17]. This implies too that the outcome of the study of [1] is the reason we agree with the effectiveness of SLMs in particular DistilBERT. Yet, according to [2], BERT or other LLMs were shown to have a major impact and did better than previous models in some areas such as disease predictions. Therefore, those models can be promising in this field. This choice of SLMs or LLMs suits the particular necessities of the application, including the effect on the computation resources, the necessity of query-time responses, as well as the importance of accuracy.

The study mentioned as one of the drawbacks the supposition of the presence of a large and diverse dataset of symptoms and their related diseases for training the machine learning models. Although in real life, sometimes the data may not be available for such a database. Furthermore, aside from measuring and comparing the inference time of the models, the survey doesn't contain the information on how these other factors could in turn affect the computation efficiency of the models, like for example the complexity of the input data or the hardware used for computation. Lastly, this study does not reflect on the ethics and privacy concerns which could be involved in operations where language models are used in health care, which are very crucial in the real world.

IX. CONCLUSION

Therefore, as the study shows both SLMs and LLMs are effective at disease prediction but there are clear allusions for improvements that guide future researchers to study behaviors of these

language models. This result shows both the options between the SLMs and LLMs and also gamuts within each mode of learning. Although LLM may be unable to improve the notion, the SLM can manage to generalize well on new data and provide fast computing speed. Therefore, SLMs provide a reliable approach in situations where the computation resources are limited or response time is fast. Nonetheless, the preference between SLMs and LLMs depends on the operational needs that range from computational resources availability to need for real time responses and predictive accuracy. Contributing to the area of forecasting disease is an important outcome of the study by elaborating the difference between the two approaches: SLMs and LLMs, furthering research, and informing the building of systems that are better and more accurate for predicting disease.

The present study explores multiple lines of research in the near future. Another area of possible research is to identify approaches that would aid in increasing the efficiency of SLMs. This could be done through approaches for model compression, quantization or hardware optimization as an example. A possibility can be to research how SLMs can be deployed in other healthcare applications, such as patient triage as well as getting medical information back. In addition, considerations of ethical and privacy issues should be included in future perspectives on SLMs for health care. It could mean looking into techniques of safeguarding personal data, obtaining written consent, and preventing the models from misusing. Finally, using a trained SLM in real life situations like a hospital or a health application and measuring their actual performance compared to others could be another exciting path for future research. This would offer innovative information about the realization and reliability of using SLMs in disease forecasting.

X. REFERENCES

[1] Sanh, V. et al. (2020) Distilbert, a distilled version of Bert: Smaller, faster, cheaper and lighter.

[2] Devlin, J. et al. (2019) Bert: Pre-training of deep bidirectional Transformers for language understanding.

[3] Schick, T. and Schütze, H. (2021) It's not just size that matters: Small language models are also few-shot learners.

[4] Aligning Large and Small Language Models via Chain-of-Thought Reasoning (Ranaldi & Freitas, EACL 2024)

[5] Brown T, Mann B, Ryder N, Subbiah M, Kaplan JD, Dhariwal P, et al. Language models are few-shot learners. *Adv Neural Inf Process Syst.* 2020; 33: 1877–901.

[6] Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, et al. Attention is all you need. *Adv Neural Inf Process Syst.* 2017; 30: 1–11.

[7] Yang R, Tan TF, Lu W, Thirunavukarasu AJ, Ting DSW, Liu N. Large language models in health care: development, applications, and challenges. *Health Care Sci.* 2023; 2: 255–263.

[8] Zhang, P. et al. (2024) TinyLlama: An open-source small language model.

[9] Hu, S. et al. (2024) MiniCPM: Unveiling the potential of small language models with Scalable Training Strategies.

[10] Awobade, B., Oduwole, M. and Kolawole, S. (2024) What happens when small is made smaller? exploring the impact of compression on small data pretrained language models.

[11]“Utilizing Natural Language Processing and Large Language Models in the Diagnosis and Prediction of Infectious Diseases: A Systematic Review.” 2024. *American Journal of Infection Control*, April.

[12] Hassan, E., Abd El-Hafeez, T. & Shams, M.Y. Optimizing classification of diseases through language model analysis of symptoms.

Sci Rep 14, 1507 (2024).

[13] StoryBuddiesPlay. 2024. *Microsoft Phi 3 LLM: Powerful Small Language Model Demystified*. StoryBuddiesPlay.

[14] Kang,Minki,et al. "Knowledge-augmented reasoning distillation for small language models in knowledge-intensive tasks." *Advances in Neural Information Processing Systems* 36 (2024).

[15] Armengol-Estapé, Jordi, et al. "SLaDe: A Portable Small Language Model Decompiler for Optimized Assembly." 2024 IEEE/ACM International Symposium on Code Generation and Optimization (CGO). IEEE, 2024.

[16] Wu, Chaoyi, et al. "PMC-LLaMA: toward building open-source language models for medicine." *Journal of the American Medical Informatics Association* (2024): ocae045.

[17] Wang, Yuling, et al. "Can Small Language Models be Good Reasoners for Sequential Recommendation?." *arXiv preprint arXiv:2403.04260* (2024).

[18] Hu, Shengding, et al. "Minicpm: Unveiling the potential of small language models with scalable training strategies." *arXiv preprint*

arXiv:2404.06395 (2024).

[19] Zhu, Yichen, et al. "LLaVA-\$\phi\$: Efficient Multi-Modal Assistant with Small Language Model." *arXiv preprint arXiv:2401.02330* (2024).

[20] Awobade, Busayo, Mardiyah Oduwale, and Steven Kolawole. "What Happens When Small Is Made Smaller? Exploring the Impact of Compression on Small Data Pretrained Language Models." *arXiv preprint arXiv:2404.04759* (2024).

[21] Magister, Lucie Charlotte, et al. "Teaching small language models to reason." *arXiv preprint arXiv:2212.08410* (2022).

[22] Fu, Yao, et al. "Specializing smaller language models towards multi-step reasoning." *International Conference on Machine Learning*. PMLR, 2023.*arXiv:2403.04260* (2024).

