# Development and evaluation of a natural language conversational bot for identifying appropriate clinician referral from patient narratives

**Sriharsha Tummala, Saptarshi Purkayastha, Ph.D., Josette Jones, Ph.D.**

Department of BioHealth Informatics, School of Informatics and Computing, Indiana University Purdue University Indianapolis

**IUPUI**
INDIANA UNIVERSITY
School of Informatics
and Computing

## Introduction

❖ Recent years have seen a significant increase in automated conversational agent chatbots. Conversational agents like chatbots for health may provide timely and cost-effective support in clinical care.

❖ Some studies show that chatbots could have an impact on patient engagement. Additionally, health systems are attempting to connect with patients over social networks, mainly where specialists are limited

❖ By 2025, the Association of American Medical Colleges estimates that the United States will have a shortfall of 61,700-94,700 physicians and critical shortage in many specialties, delaying available appointments by months in many cases.

❖ Thus, we need innovative solutions that can manage the time of limited specialists appropriately.

❖ Recent research has demonstrated that deep-learning methods are superior for natural language classification tasks compared to other machine learning methods.

❖ The primary objective of this study was to develop a telegram chatbot which reads patient narratives and acts as a conversational agent by redirecting the case to the appropriate specialist.

❖ Besides simply working on improving conversational capabilities of chatbots, we developed a novel method for referring the cases to specialists based on their responses to previous cases on a social network group.

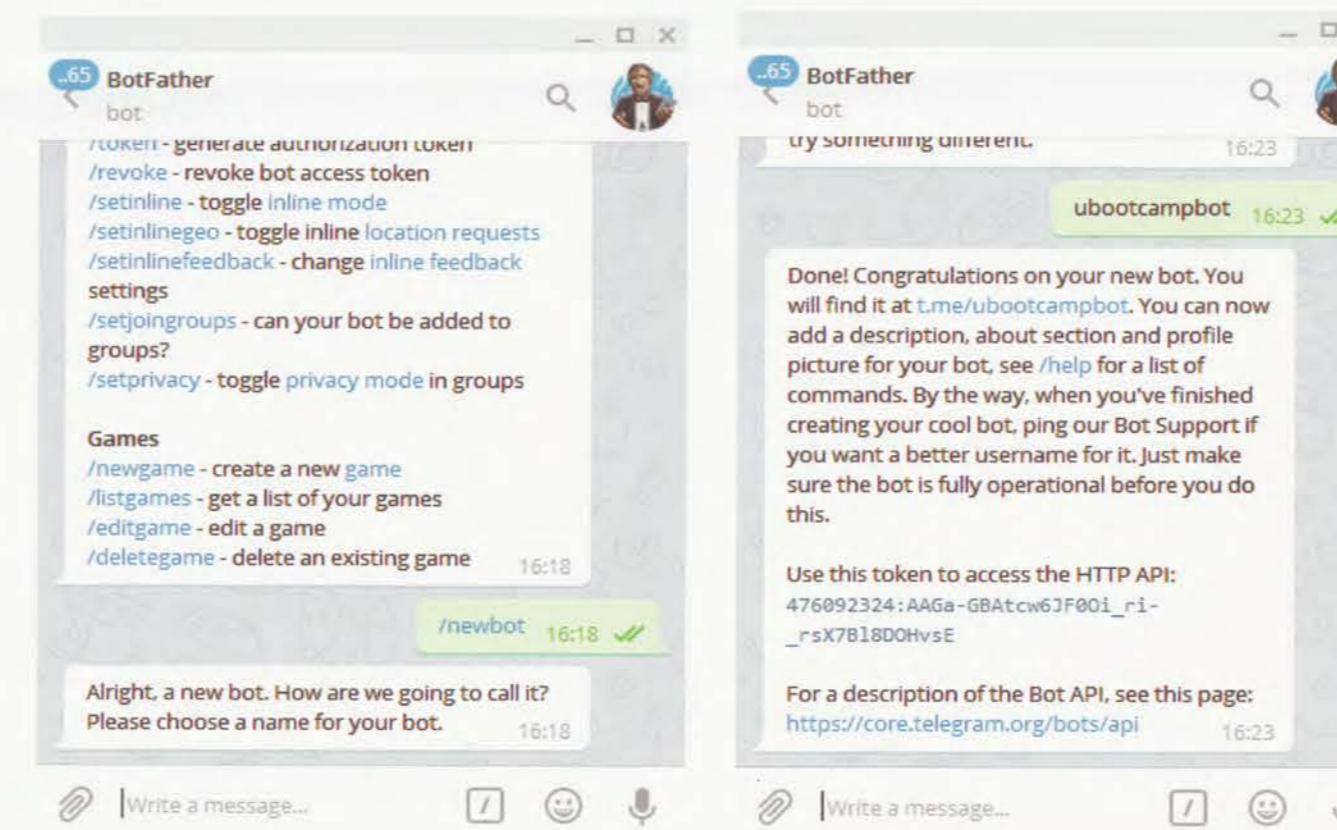❖ As far as we know, no other chatbot has the level of accuracy or referral system like our developed chatbot.

## Methodology

❖ Data is collected from Facebook consists of 1890 clinicians. The data is 2 years old and it is deidentified

❖ 568 actors identified

❖ Inclusion Criteria: Must have made at least 1 post or comment

❖ 1,143 (top level) posts and 8,606 comments made to these posts

❖ *Commauth* and *Postcontent* are extracted as primary columns from the data set and used for model development. Data was 3x imputed to get appropriate size for training a deep neural network.

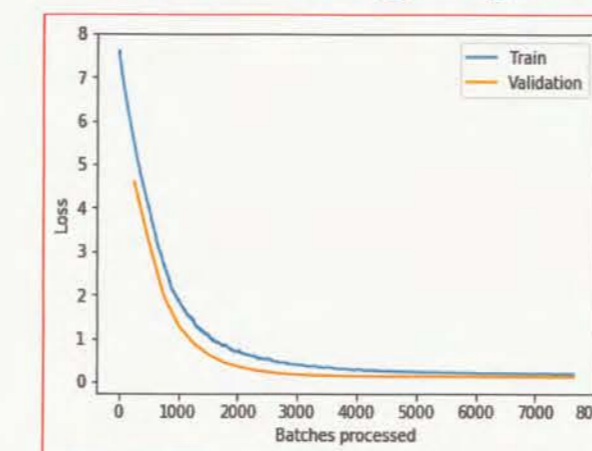| | label | text |
|---|---|---|
| 0 | boudhayan.dm | Dr MpSingh, Dr Swagatawhat would you suggest f... |
| 1 | durga.prasan | Dr MpSingh, Dr Swagatawhat would you suggest f... |
| 2 | durga.prasan | Ideally the rectal polyp should have been snar... |
| 3 | durga.prasan | My suggestion is that all cancers should be de... |
| 4 | swagata.brahmachari | Ya surely .As correctly pointed out by Durga P... |

❖ Deep learning library *fastai* is used for the model development and training.

❖ Tokenizer was used to tokenize the important vocabulary. Pretrained weights from Wikipedia trained model was downloaded as part of transfer learning.

❖ Data was divided into training and target using TextDataBunch

❖ The language model was trained for 30 epochs, which resulted in an accuracy of 97% with a converging training (0.190) and validation loss (0.119).

❖ The language model creation process resulted in a neural network with an embedding size of 400, 3 layers, 1150 hidden activations per layer.
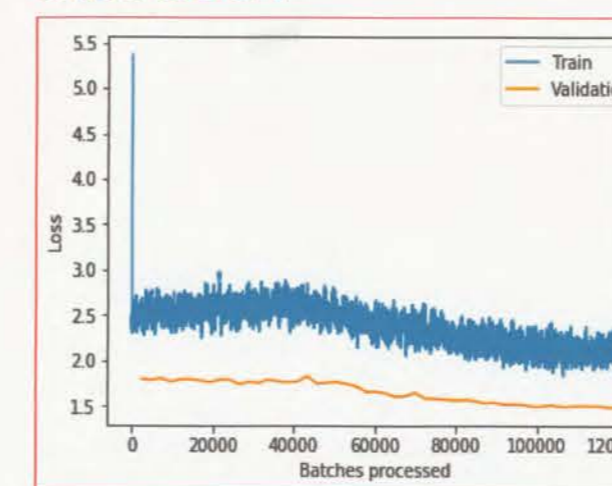
## Software developed
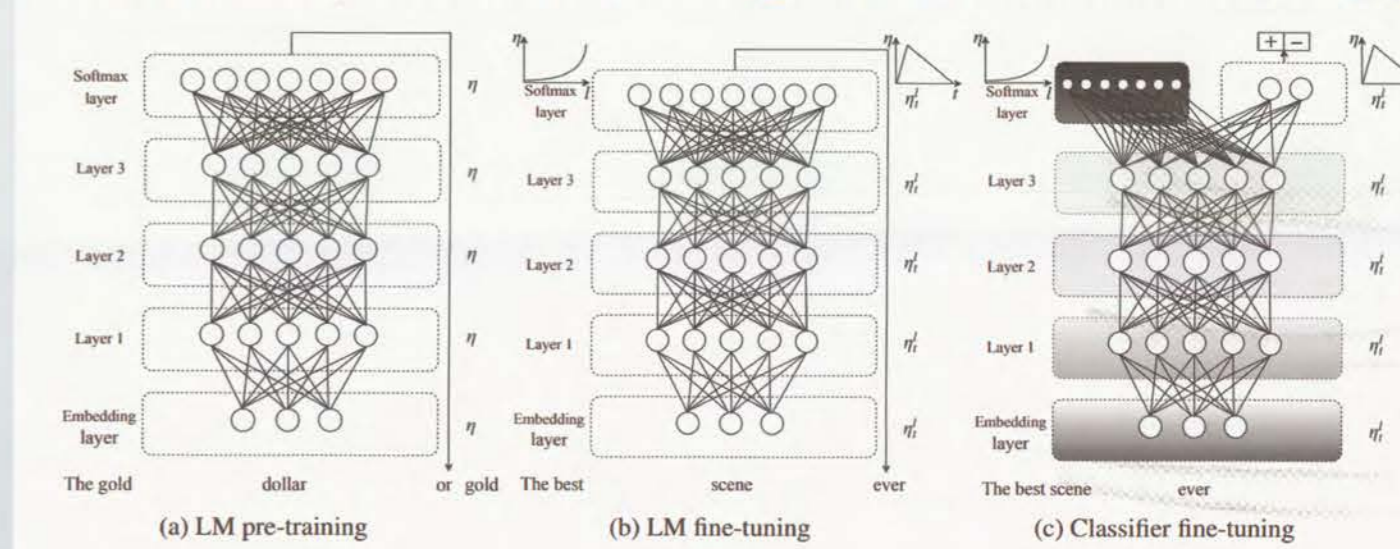


## Results

❖ Losses in language model at 30 epochs



❖ The pre-trained vocabulary from Wikipedia is used for transfer learning for the NLP using the ULMFiT (Universal Language Model Fine-tuning for Text Classification) [1]

❖ We then trained a classified using the AWD-LSTM algorithm with a batch size of 16 and with 3 different learning rates from 0.01 to 0.000001



## The ULMFiT transfer learning for NLP classification [1]



(a) LM pre-training  (b) LM fine-tuning  (c) Classifier fine-tuning

❖ The classifications have a 59% accuracy, but they are close to 100% accurate based on gold-standard human verification.

❖ Since there are more specialists of the same specialization in the social network, the classifier approximates the specialists, which was then verified manually to be accurate.

## Limitations

❖ Data collected was very low in size to be suitable for Neural networks

❖ The data consists many of non-English words, images in the posts and typos which poses difficulty in training and developing the model

## Future work and conclusion

❖ To deploy this on patient portals or a Facebook/telegram/WhatsApp group and evaluate the clinician response, and maybe later patient outcomes.

❖ We believe further research of this approach can help reduce burden, and ease the load of resource limited health systems.

References:

[1] Howard, Jeremy, and Sebastian Ruder. "Universal language model fine-tuning for text classification." In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), vol. 1, pp. 328-339. 2018.

[2] H. Wang, Q. Zhang, M. Ip and J. T. Fai Lau, "Social Media–based Conversational Agents for Health Management and Interventions," in Computer, vol. 51, no. 8, pp. 26-33, 2018. doi:10.1109/MC.2018.3191249