

Panel: Methods and Applications of Biomedical Natural Language Processing Across Languages and Institutions

Hua Xu, Rui Zhang, Yonghui Wu, Buzhou Tang, Gayo Diallo





Agenda

- Introduction (5 min)
- Presentation
 - Each of 5 panellists present 4 min (20 min)
- Discussion (15 min)



Panellists



Hua Xu, PhD,
FACMI
Professor
Yale University
USA



Rui Zhang, PhD
Associate professor
Univ of Minnesota
USA



Yonghui Wu, PhD
Associate Professor
University of Florida
USA



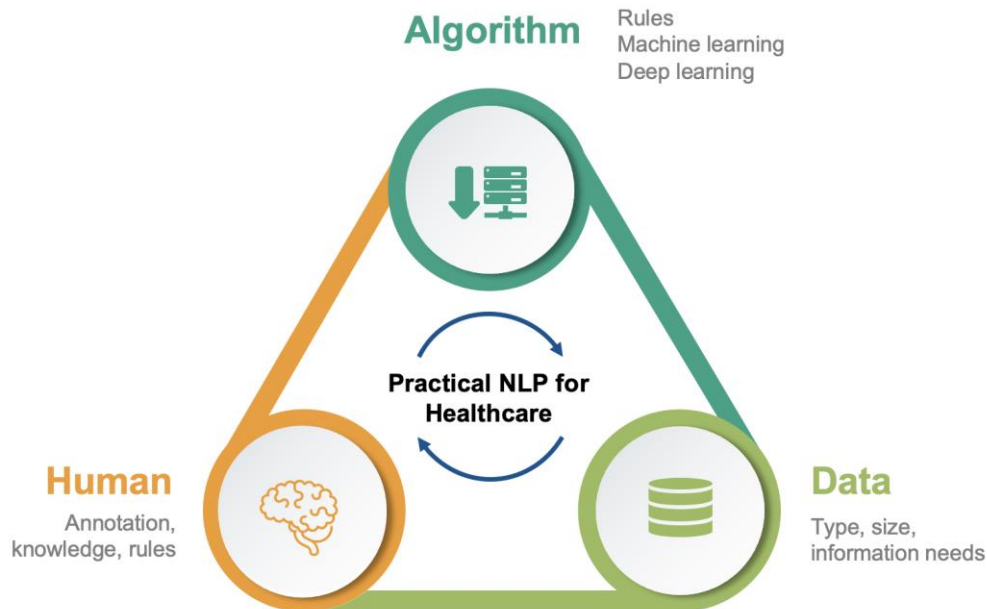
Buzhou Tang, PhD
Associate Professor
Herbin Institute of
Engineering
China



Gayo Diallo
Professor
Univ. of Bordeaux
France

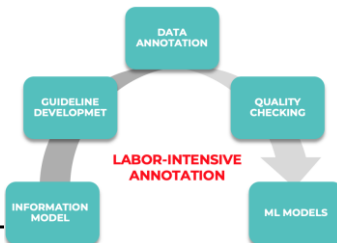
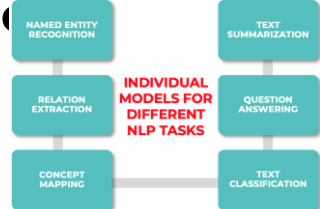


Current clinical NLP development process





How LLMs will change current practice?



- Phlebotomist
- Respiratory therapist
- Blood gas
- Needle
- Protective clothing
- Gloves
- Arm
- Index finger
- Emergency department
- Baseline studies
- Employee health
- HIV
- Hepatitis C
- Periodic screening
- Blood tests
- Final exam

SHOT

HU

Extract without rephrasing all treatment entities from the following note in a list format: "HISTORY OF PRESENT ILLNESS : The patient is a 54 - year - old right - handed male who works as a phlebotomist and respiratory therapist at Hospital . The patient states that he was attempting to do a blood gas . He had his finger of the left hand over the pulse and was inserting a needle using the right hand . He did have a protective clothing including use of gloves at the time of the incident . As he advanced the needle , the patient jerked away , this caused him to pull out of the arm and inadvertently pricked the tip of his index finger . The patient was seen and evaluated at the emergency department at the time of incident and had baseline studies drawn , and has been followed by employee health for his injury . The source patient was tested for signs of disease and was found to be negative for HIV , but was found to be a carrier for hepatitis C . The patient has had periodic screening including a blood tests and returns now for his final exam . "

Glucophage	850	bid
Glipizide	10	bid
Imodium	Not specified	prn

Zero-shot Clinical Entity Recognition using ChatGPT

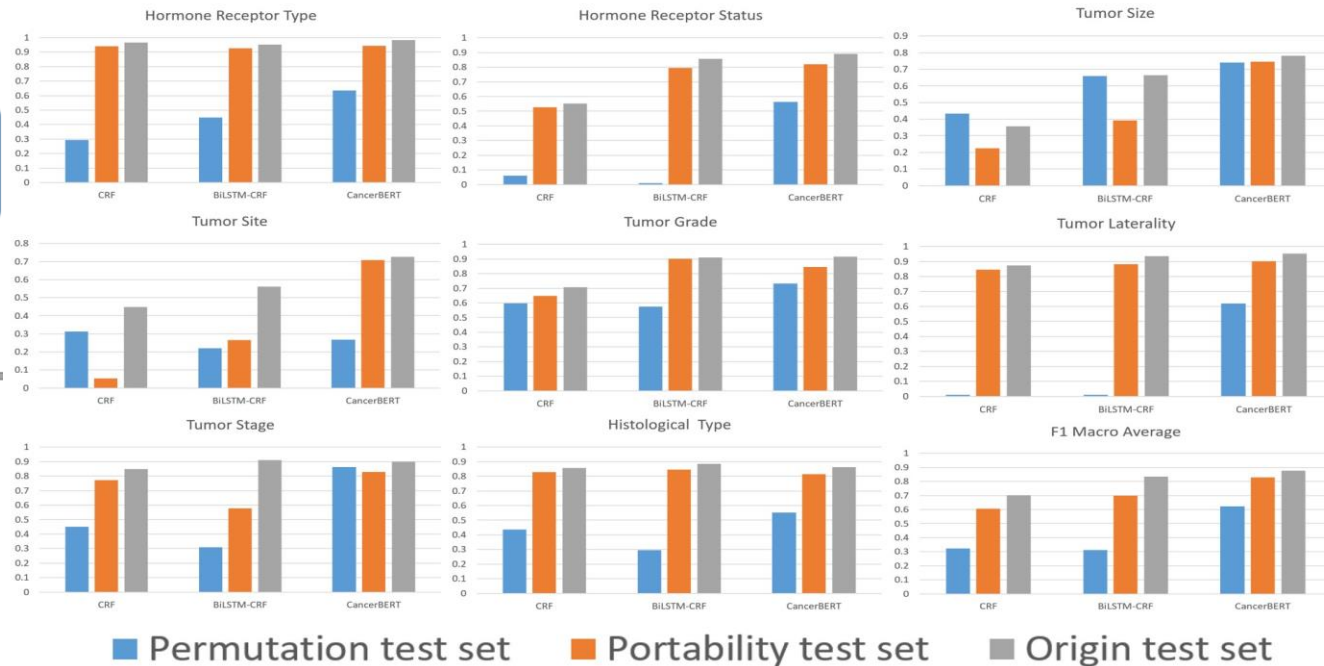
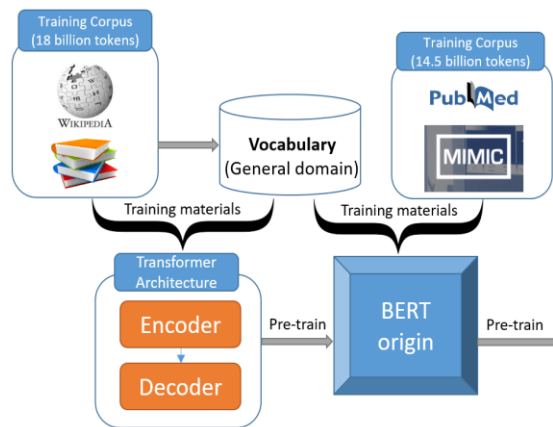
Yan Hu, Iqra Ameer, Xu Zuo, Xueqing Peng, Yujia Zhou, Zehan Li, Yiming Li, Jianfu Li, Xiaoqian Jiang, Hua Xu

In this study, we investigated the potential of ChatGPT, a large language model developed by OpenAI, for the clinical named entity recognition task defined in the 2010 i2b2 challenge, in a zero-shot setting with two different prompt strategies. We compared its performance with GPT-3 in a similar zero-shot setting, as well as a fine-tuned BioClinicalBERT model using a set of synthetic clinical notes from MTSamples. Our findings revealed that ChatGPT outperformed GPT-3 in the zero-shot setting, with F1 scores of 0.418 (vs.0.250) and 0.620 (vs. 0.480) for exact- and relaxed-matching, respectively. Moreover, prompts affected ChatGPT's performance greatly, with relaxed-matching F1 scores of 0.628 vs.0.541 for two different prompt strategies. Although ChatGPT's performance was still lower than that of the supervised BioClinicalBERT model (i.e., relaxed-matching F1 scores of 0.628 vs. 0.870), our study demonstrates the great potential of ChatGPT for clinical NER tasks in a zero-shot setting, which is much more appealing as it does not require any annotation.

Research and Applications

CancerBERT: a cancer domain-specific language model for extracting breast cancer phenotypes from electronic health records

CancerBERT



4.5 M clinical notes; 1.3M pathology reports; 1B tokens

NLP for drug repurposing

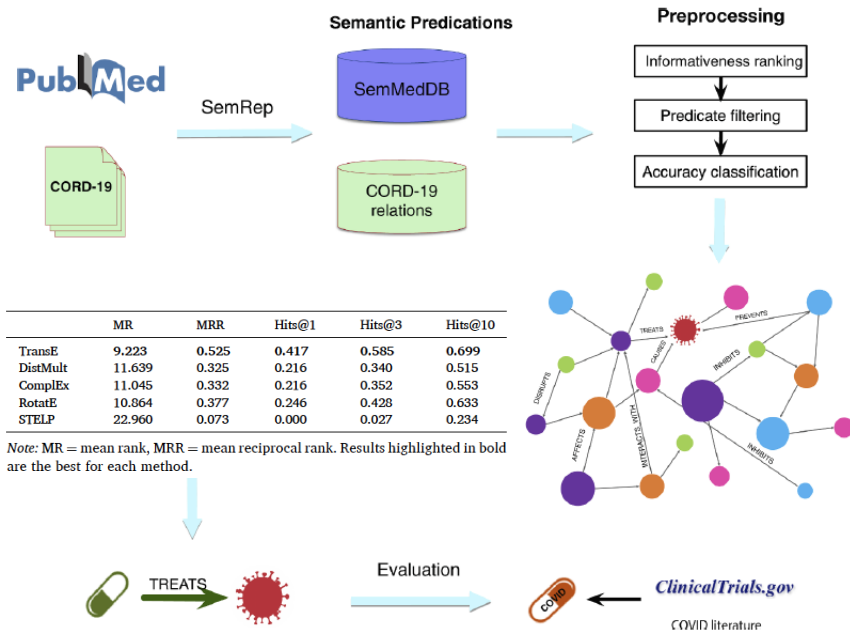
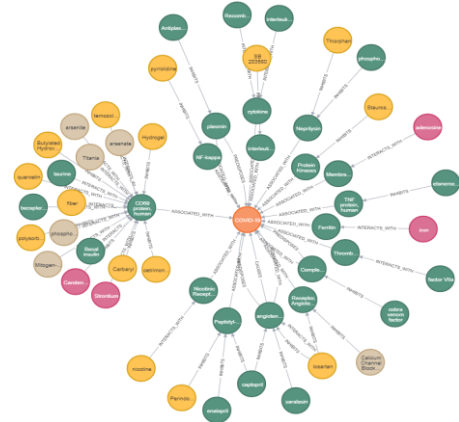
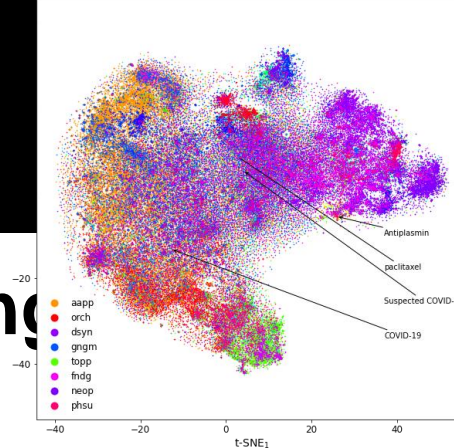


Fig. 1. Diagram illustrating the workflow of our approach.

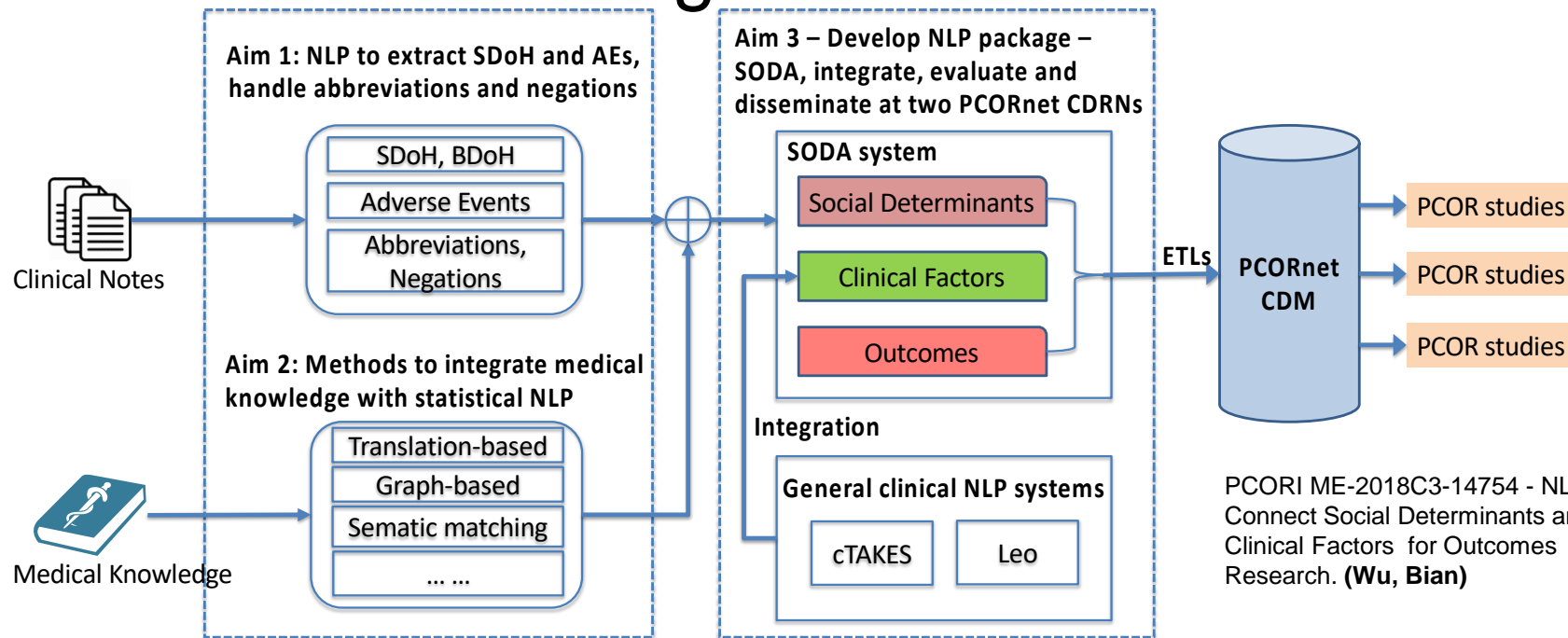


Thirty-three candidate drugs highly ranked by TransE and deemed plausible in manual analysis.

Metoclopramide	Trilostane
Oxymatrine	Cyproterone Acetate
Mitogen-Activated-Protein Kinase Inhibitor	Nucleoside Reverse-Transcriptase Inhibitors
Oxophenylarsine	Methyltrienolone
5-Alpha reductase inhibitor	Bosentan
Folic acid	Estramustine
Anthelmintics	Allicin
Sildenafil	Proteasome inhibitors
Furosemide	Antiplatelet Agents
Beclomethasone	Fibrinolytic Agents
Cangrelor	Contraceptive Agents
Gymnemic acid	Neuraminidase inhibitor
Estradiol	Vitamin D Analogue
mTOR Inhibitor	Tyrosine kinase inhibitor
Clobetasol propionate	Mometasone furoate
Carbenoxolone	Vasopressin Antagonist
Anti-Retroviral Agents	



Extract SDoH Using NLP



PCORI ME-2018C3-14754 - NLP to Connect Social Determinants and Clinical Factors for Outcomes Research. (Wu, Bian)

GatorTron models : 345 Million, 3.9 Billion, and 8.9 Billion parameters

npj | digital medicine

www.nature.com/npjdigitalmed

ARTICLE OPEN

Check for updates

A large language model for electronic health records

Xi Yang^{1,2}, Aokun Chen^{1,2}, Nima PourNejatian³, Hoo Chang Shin³, Kaleb E. Smith³, Christopher Parisien³, Colin Compas³, Cheryl Martin³, Anthony B. Costa³, Mona G. Flores³, Ying Zhang³, Tanja Magoc³, Christopher A. Harle^{1,5}, Gloria Lipori^{5,6}, Duane A. Mitchell⁶, William R. Hogan¹, Elizabeth A. Shenkman¹, Jiang Bian^{1,2} and Yonghui Wu^{1,2,6}

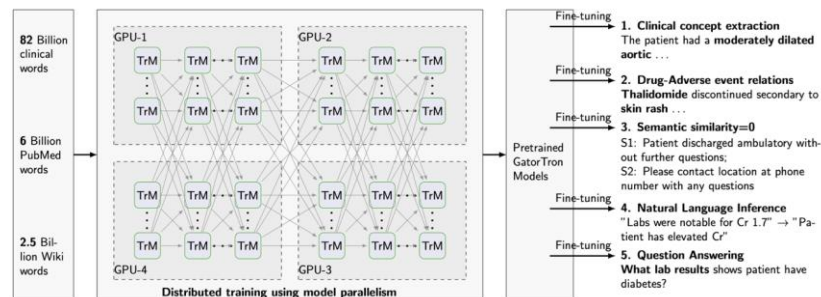
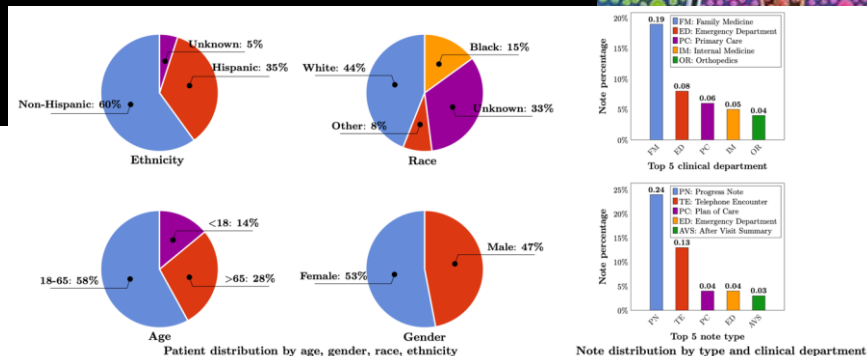
There is an increasing interest in developing artificial intelligence (AI) systems to process and interpret electronic health records (EHRs). Natural language processing (NLP) powered by pretrained language models is the key technology for medical AI systems utilizing clinical narratives. However, there are few clinical language models, the largest of which trained in the clinical domain is comparatively small at 110 million parameters (compared with billions of parameters in the general domain). It is not clear how large clinical language models with billions of parameters can help medical AI systems utilize unstructured EHRs. In this study, we develop from scratch a large clinical language model—GatorTron—using >90 billion words of text (including >82 billion words of de-identified clinical text) and systematically evaluate it on five clinical NLP tasks including clinical concept extraction, medical relation extraction, semantic textual similarity, natural language inference (NLI), and medical question answering (MQA). We examine how (1) scaling up the number of parameters and (2) scaling up the size of the training data could benefit these NLP tasks. GatorTron models scale up the clinical language model from 110 million to 8.9 billion parameters and improve five clinical NLP tasks (e.g., 9.6% and 95% improvement in accuracy for NLI and MQA), which can be applied to medical AI systems to improve healthcare delivery. The GatorTron models are publicly available at: https://catalog.ngc.nvidia.com/orgs/nvidia/teams/clara/models/gatortron_og.

npj Digital Medicine (2022)5:194; <https://doi.org/10.1038/s41746-022-00742-2>

Models available from: <https://huggingface.co/UFNLP>

@TheInstituteDH #MEDINFO23

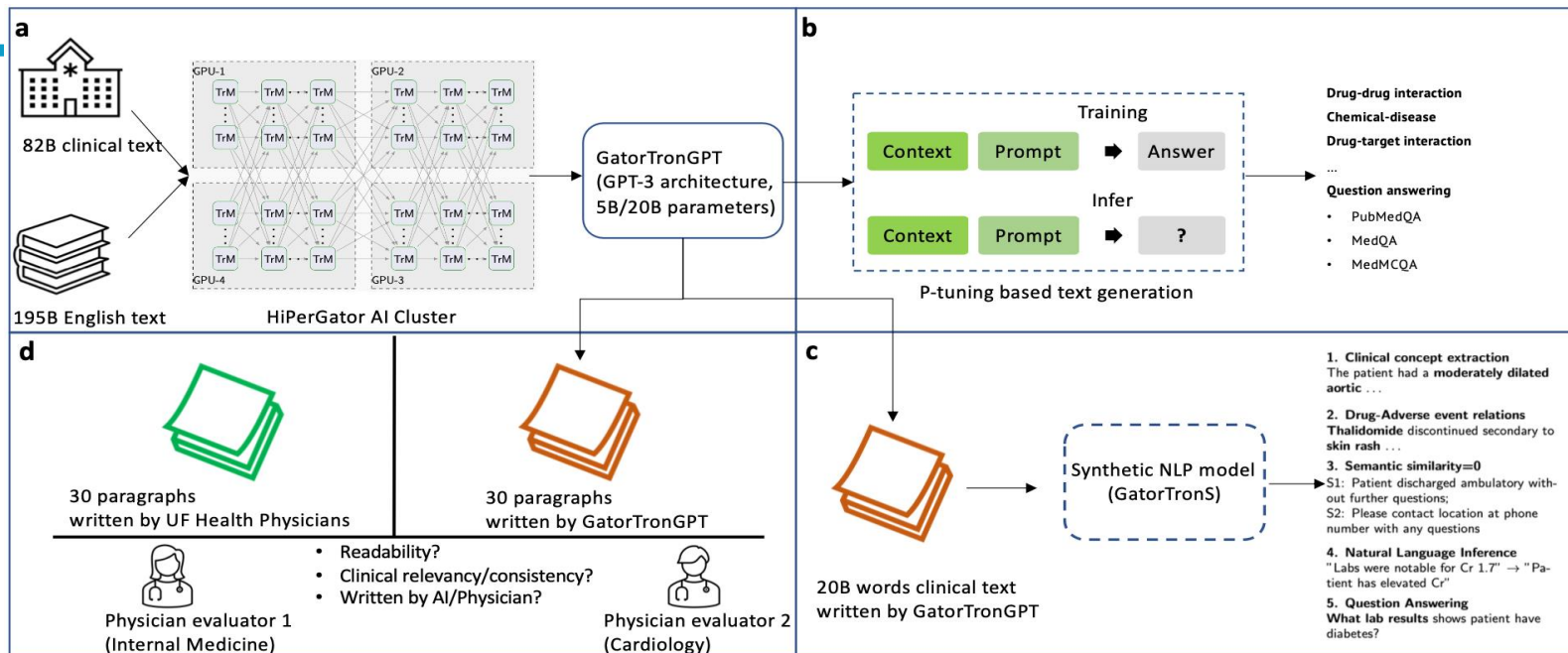
Yang X, Chen A, PourNejatian N, Shin HC, Smith KE, Parisien C, Compas C, Martin C, Costa AB, Flores MG, Zhang Y, Magoc T, Harle CA, Lipori G, Mitchell DA, Hogan WR, Shenkman EA, Bian J, Wu Y. **A large language model for electronic health records.** Npj Digit Med. Nature Publishing Group; 2022 Dec 26;5(1):1–9.



Model	# Layers	# Hidden Size	# Attention Heads	# Parameters
GatorTron-base	24	1024	16	345 million
GatorTron-medium	48	2560	40	3.9 billion
GatorTron-large	56	3584	56	8.9 billion



GatorTronGPT - A generative LLM for EHRs



Phenotyping and automatic coding

- ① More interactions, more gains
- ② Existing knowledge graphs is beneficial to many tasks including phenotyping and automatic coding

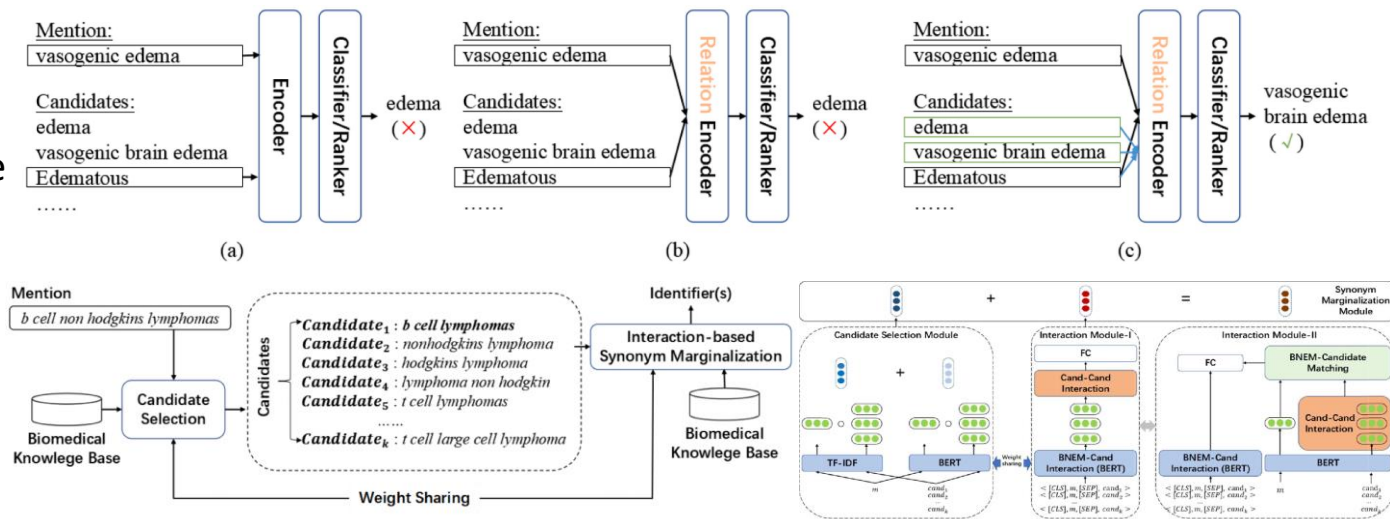
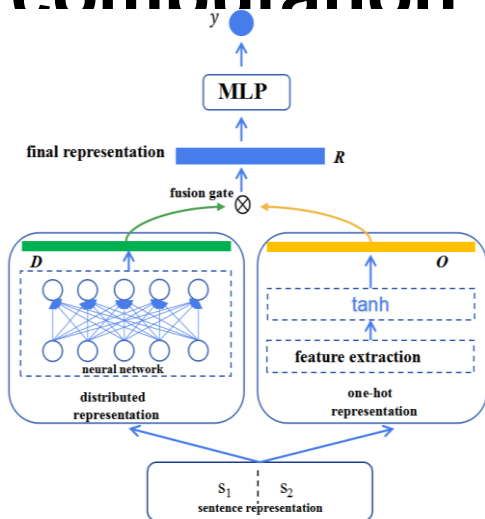


Fig. 2. Overview of our method (IA-BIOSYN).

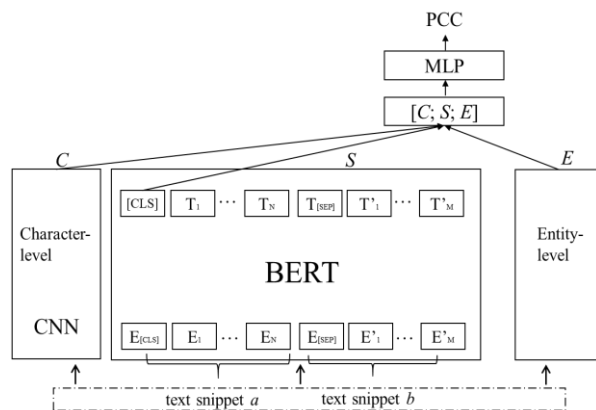
The architecture of IA-BIOSYN

Peng H, Xiong Y, Xiang Y, Wang H, Xu H, Tang B. Biomedical named entity normalization via interaction-based synonym marginalization. J Biomed Inform. 2022 Dec;136:104238. doi: 10.1016/j.jbi.2022.104238. Epub 2022 Nov 15. PMID: 36400329.

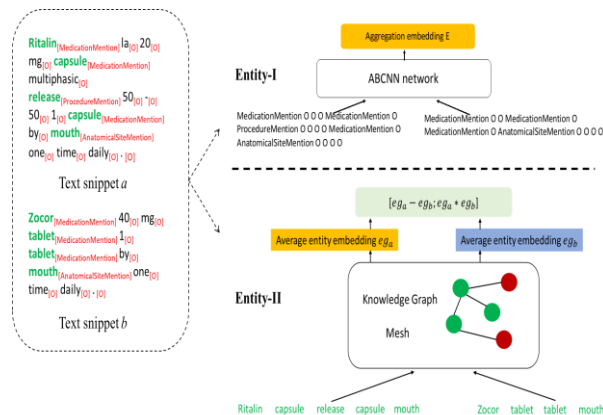
Clinical semantic textural similarity computation



① distributed representation
+ one-hot representation



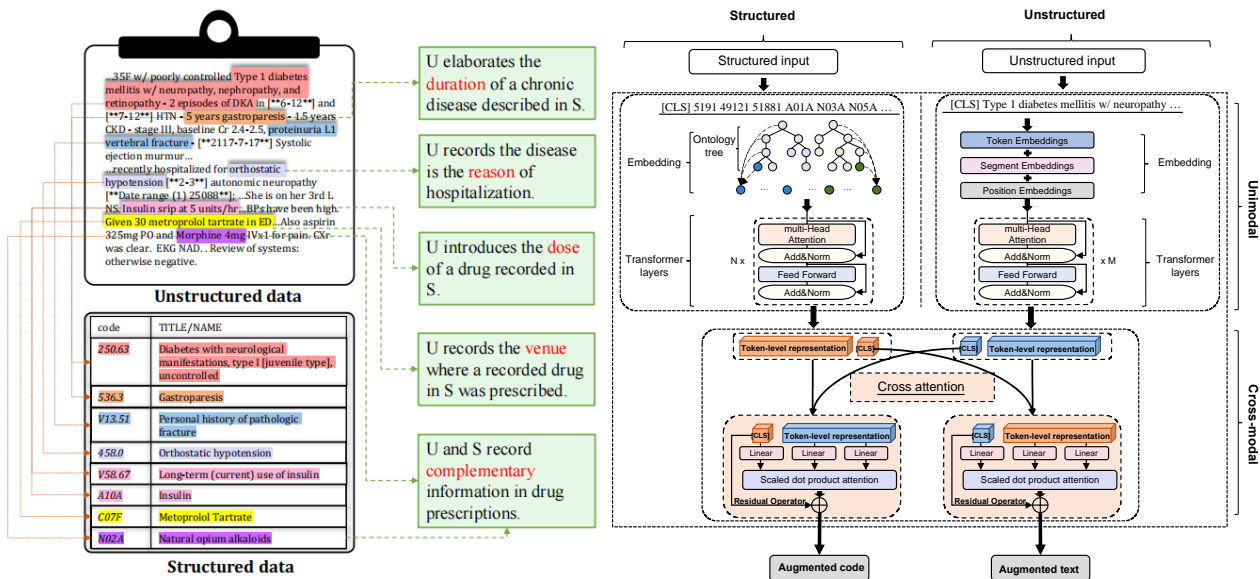
② multi-granularity similarity
fusion



③ soft-alignment for
similarity computation

Multimodal language models and clinical data analysis

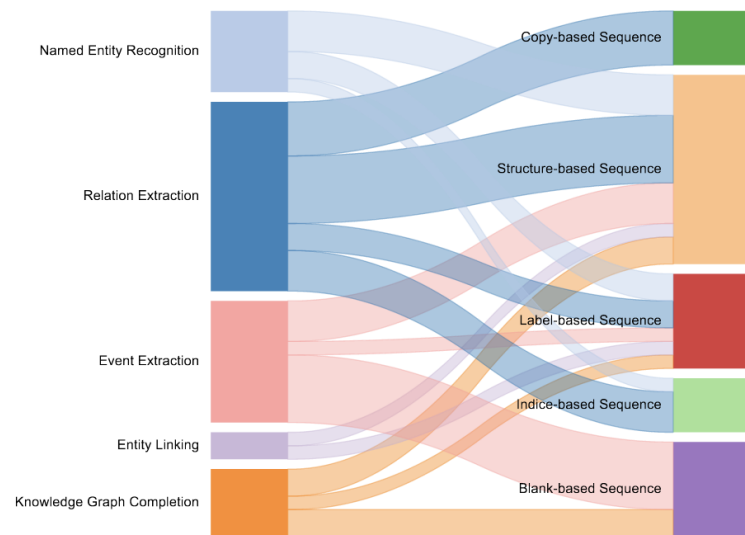
- There are relations between structured data and text.
- Design two pretraining tasks: Text2Code, Code2Code.
- Multimodal language models are better than monomodal models on NLP and outcome prediction tasks.





NLP & BioMed Knowledge Graphs acquisition

- A Knowledge Graph
 - Is a machine-readable representation of domain-specific knowledge
 - Is a directed labeled graph with clearly defined labels.
 - Components are nodes, edges, and labels
- NLP & LLMs for designing and population
 - Named Entity Recognition
 - Relation Extraction
 - Event Extraction
 - Entity Linking
 - Multilingual labels acquisition for non-English languages



Source: Ye et al. (2022): Generative Knowledge Graph Construction: A Review
EMNLP 2022

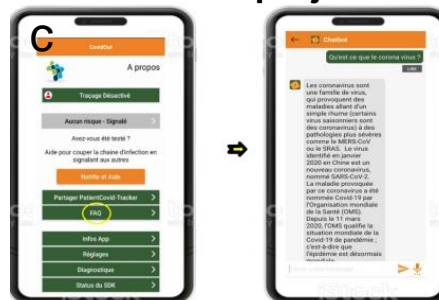


How NLP could contribute to more inclusive Digital Health in Low & Middle Income Countries

› Health Literacy and eHealth Literacy are key issues

- Taking into account under-resourced languages
- NLP beyond written texts
- Producing annotated corpora
- Voice processing technology for Text-To-Speech and Speech-to-Text

The CovidBot of the PATIENT-COVID19 project



Improving Covid-19 vaccine literacy among undergraduate students in Burkina Faso

Michel J Some; Ismaila Ouedraogo; Roland Benedikter; Rasmané Yameogo; Ghislain Ateamezing; Ibrahim Traoré; Gayo Diallo
2022 IEEE 10th International Conference on Healthcare Informatics (ICHI)
Year: 2022 | Conference Paper | Publisher: IEEE

Abstract HTML PDF

The Covid-19 pandemic has had a major effect on education. University students are going through a high level of psychological pressure and the pandemic has called for people to seek and use covid19 related information to adapt their behaviours. In this study, we present an AI-enabled **Interactive Voice Response Service to Improve High School Students Covid-19 Literacy in Burkina Faso: A Usability Study** conduct an impact study

Authors Michel J. Some, Ismaila Ouedraogo, Roland Benedikter, Rasmané Yameogo, Ghislain Ateamezing, Ibrahim Traoré, Gayo Diallo
Pages 454 - 457
DOI 10.3233/SHTI220763
Category Research Article
Series Studies in Health Technology and Informatics
Ebook Volume 295: Advances in Informatics, Management and Technology in Healthcare

Abstract

Mobile technology is widely used in healthcare. However, designers and developers in many cases have focused on developing solutions that are often tailored to highly literate people. While the advent of the pandemic has called for people to seek and use Covid-19 related information to adapt their behaviors, it is relatively difficult for low literate to get easily access to health information through digital technologies. In this study, we present a Mobile based Interactive Voice Response service designed particularly for low-literate people which provides validated Covid-19 related health information in local African languages. We conducted a field study, among high school students, through a usability study to assess users' perception. The service received an excellent numerical usability score of 78.75

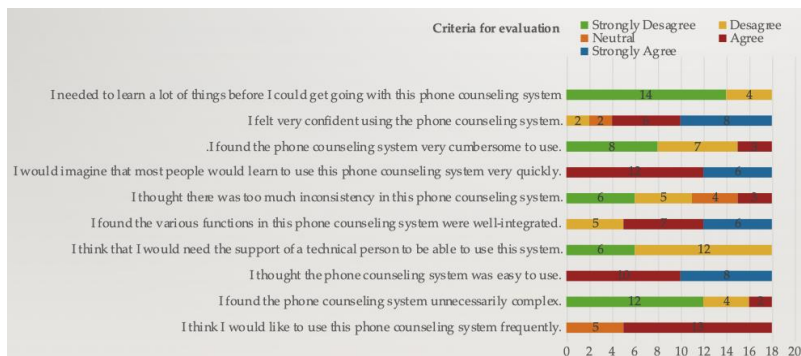


How NLP could contribute to more inclusive Digital Health in Low & Middle Income Countries

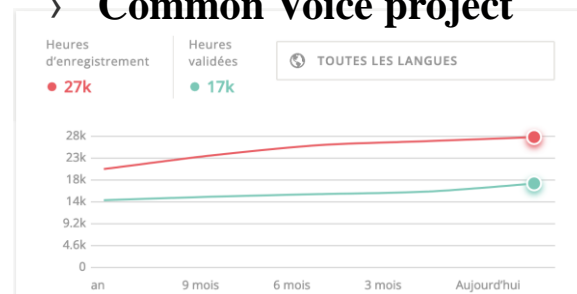


Evaluation of CovidtBot: System Usability Scale

Participants (n=18) low literates (without high diploma) , 12 were women (66.7%) and 6 were men (33.3%). 12 < Age < 20 years old.



Contribution to the Mozilla Common Voice project



Français

Heures: 1076
Locuteurs et locutrices: 17778
Avancement de la validation: 89%

CONTRIBUER

Dioula

Heures: 1
Locuteurs et locutrices: 1
Avancement de la validation: 100%

KA BÓLOMAFARA DI



Panel Discussion

- Challenges in NLP methodology and clinical practical applications
- Opportunities for NLP in clinical research and practice
- Challenges of clinical NLP at different linguistic settings and different countries
- Challenges and experience of applying deep learning and/or large language models to real-world