CardioNeoCanon: a customized chatbot for cardiology inquiry with retrieval augmented generation to reduce hallucinations and improve performance of large language models

Tu Hao Tran **CSANZ**

August 2nd 2024

Faculty Disclosure

The presenter has advised that the following presentation will NOT include discussion on any commercial products or service and that there are NO financial interests or relationships with any of the Commercial Supporters of this years ASM.

Aims, Background and Method

- Large Language Models (LLMs) and chatbots are becoming ubiquitous.
- Adoption in cardiology practice is understandably hampered by misleading information and/or hallucinations (fiction).
- A customised chatbot for use by cardiologists needs to be:
 - Accurate
 - Not misleading
 - Up-to-date
 - Robust referencing.



CardioNeoCanon compared with **ChatGPT 3.5** for 4 questions by 11 **blinded** cardiology doctors:

- 1. Create a detailed summary of an approach to the management of heart failure with evidence-based practice.
- 2. List the different types of CIED and give the indications, follow-up period in months and clinical considerations, in a table.
- 3. Compare paragon-HF and paradigm-HF trials in table
- 4. Propose a framework for a pHD which aims to elucidate the pathophysiological substrate for atrial fibrillation that involves rotors, artificial intelligence and augmented reality. Give a detailed summary of the field of rotors in atrial fibrillation and give 5 high impact articles that will assist with this field of study.



- All doctors subjectively adjudicate the performance with a **VRISM** score (5-15, low=good, high=poor):
 - Vagueness (1 not vague, 2 vague, 3 very vague)
 - Referencing (1 references all correct, 2 no references, 3 references incorrect)
 - Incorrect information (1 none, 2 not sure, 3 incorrect)
 - **S**tructure (1 good, 2 adequate, 3 poor)
 - **M**issing critical information (1 none, 2 not sure, 3 missing)
- One-tailed Wilcoxon signed-rank test:
 - Small sample size without normal distribution
 - W-value of 1 with a mean difference of 8.73.
 - Indicating that CardioNeoCanon achieved significantly lower VRISM scores compared to ChatGPT 3.5 (p<0.05)

Treatment 1 23 27 22 31 32 22 27 35 35 35 31 26	Treatment 2 39 37 34 36 35 34 41 42 33 41 39	Sign -1 -1 -1 -1 -1 -1 -1 -1 -1 -1	Abs	16 10 12 5 3 12 14 7 2 10 13	R 11 5.5 7.5 3 2 7.5 10 4 1 5.5 9	Sign R -11 -5.5 -7.5 -3 -2 -7.5 -10 -4 1 -5.5 -9	
Significance Level:					Result Details W-value: 1		
•.05					Mean Difference: -8.73 Sum of pos. ranks: 1 Sum of neg. ranks: 65		
1 or 2-tailed hypothesis?:					<i>Z</i> -value: -2.8451 Mean (<i>W</i>): 33		
 One-tailed 					Standard Deviation (<i>W</i>): 11.25		
⊖ Two-tailed					Samp	le Size (<i>N</i>): 11	

Conclusion and Discussion

- **CardioNeoCanon** achieved significantly lower VRISM scores by a considerable margin (p<0.05)
- Demonstrates the practicality of clinicians developing their own chatbots that can **outperform** off the shelf LLM options.
- Clinicians should strive to have more **agency** in the space of artificial intelligence and embrace participating in **innovation**.

