Do Currently Available Large Language Models Provide Musculoskeletal Treatment Recommendations That Are Concordant with Evidence-Based Clinical Practice Guidelines?

Benedict U Nwachukwu¹, Nathan Varady, Answorth Anthony Allen², David W Altchek², Joshua S Dines¹, Riley Joseph Williams³, Kyle N Kunze⁴

¹Hospital For Special Surgery, ²Hosp for Special Surgery, ³Hospital For Special Surgery - Weill Cornell Med, ⁴Hospital for Special Surgery

INTRODUCTION: An important application of large language models (LLMs) consists of the feasibility for providing evidenced-based guidelines for medical providers, as this may help guide patient management including injury triage and necessity of referral to musculoskeletal specialists, and therefore have implications for deployment within healthcare settings. The purpose of this study was to determine whether several leading, commercially-available LLMs provide treatment recommendations concordant with evidenced-based clinical practice guidelines (CPGs) developed by the American Academy of Orthopedic Surgeons (AAOS).

METHODS: All evidence-based CPGs concerning the management of rotator cuff tears (n=33) and anterior cruciate ligament (ACL) injuries (n=15) were extracted from the AAOS (**Table 1**). Utilizing a structured framework for query development, information from each AAOS guideline was transformed into a question-based format to prompt each LLM to provide a recommendation concerning the topic. Subsequently, treatment recommendations from four contemporary LLMs (Chat-generative pretrained transformer version-4 [ChatGPT-4; OpenAI], Gemini (Google), Mistral-7B (Mistral AI), and Claude-3 (Anthropic) were obtained. Each question was systematically queried three times utilizing each LLM, and additional follow-up dialogue was not permitted after the initial query. After the three prompts, responses were categorized into a single unique recommendation for analysis by two blinded physicians based on majority vote (576 total responses were therefore combined into 192 unique recommendations for analysis) as being "concordant," "discordant," or "indeterminate" (i.e., neutral response without definitive recommendation) with respect to AAOS CPGs. Prior to each new query, the previous dialogue was deleted and history cleared in order to avoid the propagation of bias from stored memory stored. The overall concordance between LLM and AAOS recommendations were quantified, while the comparative overall concordance of recommendations amongst the four LLMs was evaluated through the Fischer's-exact test. Inter-rater reliability of response concordance was assessed utilizing Cohen's Kappa coefficient. RESULTS:

A total of 192 responses were elicited. Overall, 135 (70.3%) responses were concordant, 43 (22.4%) were indeterminate, and 14 (7.3%) were discordant (**Table 2**). Inter-rater reliability for classification of concordance was deemed as excellent (Kappa=0.92). Concordance with AAOS CPGs was most frequently observed with ChatGPT-4 (n=38, 79.2%), and least frequently with Mistral-7B (n=28,58.3%). Indeterminate recommendations were most frequently observed with Mistral-7B (n=17,35.4%) and least frequently with Claude-3 (n=8, 6.7%). Discordant recommendations were most frequently observed with Gemini (n=6,12.5%) and least frequently with ChatGPT-4 (n=1,2.1%). Overall, no statistically significant differences in concordant recommendations was observed across LLMs (p=0.12; **Figure 1**). Only 20 (10.4%) of all recommendations were transparent and provided references with full bibliographic details or links to specific peer-reviewed content to support recommendations (**Figure 2**) ChatGPT-4 was the most transparent LLM evaluated, providing 11 (22.9%) responses with complete and transparent citations. Mistral-7B was the least transparent LLM, with 45 (93.4%) of responses not providing any reference or source to support the recommendation.

DISCUSSION AND CONCLUSION: More than one-in-four recommendations provided by leading commercially-available LLMs concerning the evaluation and management of rotator cuff and ACL injuries are not concordant with current evidenced-based CPGs. Although ChatGPT-4 demonstrated the highest performance, clinically significant rates of recommendations without concordance or supporting evidence exist across all the evaluated LLMs. Only 10% of responses by LLMs were transparent and provided complete, evidence-based resources, precluding users from fully understanding and interpreting the sources from which recommendations were provided. Future academic research should include analyses of multiple LLMs beyond ChatGPT given that performance may vary; in the interim, while all leading LLMs generally provide recommendations concordant with CPGs, the substantial proportion of recommendations that do not align with these guidelines suggest that LLMs are not adequate clinical support tools.





Kabgori	AADS Recommendation	Investigation of Concernment of Conc	
Management of Small to Medium Publishings Team*	Broug evidence supports that both physical through and operative transmit work to cipalitant importance is puttern exported success for patients with componentic smallers making fail- functions there of them.	Georg	
Management of Figh Goale Partial Distance Team*	Being existence apports for use of either sometrismen hill disistence or transmissions in site report is potents that failed conservative management with high-goale partial disistence transmi- crift test.	Georg	
Long term Nan apendra Mangamont ^a	Sense evidence appoint that palent reported nationes: improve with physical theory in components: partness with full fieldness rotates call tasts. However, the neural call tast size, much structure, and tably infitiation may program over 5 to 30 pices with non-operative measurement.	Georg	
Bigooli (Chied Exmination)*	Strag endoge apports for dising continuing on he such to dispose or straft prices with opper cell trans, however, continuing of note will move disposite access.	Story	
Elaposis (Inaging)*	String evidence apparts for XBE, VRA, and altraceast are welated parts to a clinical exam- for identifying exister cuff ture.	Story	
Past operative Mobilization Taning*	Sense evidence acquires similar paraperative efficient and patient expands managers for small to random aired that disclosure renter and more between only notification and defauld model taking in a Y work to patient which are undergone attheought, interest of require	Soung	
Proposite Factors (Apr)*	foring evidence apports that dder up is associated with higher failure rates and poerty parient reported entropies whet totate coll repair.	Storeg	
Pergnandic Factors (Worker's Communication)*	Brong evidence apports the preserve of a series's comparation claim is associated with more ration reported surveyses after strates cell series	Georg	
Biologic Approximition with Fatcher Destroid Products?	Broug evidence due not support biological supported internet, foreire cell waits with platele- dering produces an importing patient reported internets, however, limited evidence support. Her use of liquid platelet relation plasma (PP) is the context of descending to our root.	lines	
Engle Ken vs. Chaitie Ken Kepuir ⁴	Brong evidence doe not support double way rotator cull upon constructs on improving nation constructionscenary compared to shark you carried nationa trust constructs.	Lincog	
lingle Ken vs. Dealth Ken Rapain. In Re Tour?	Strong evidence supports lower re-true usins after dealth ever supin compared to single rev- variant matters repair where evaluating for both partial and full features strates after primary repair. Network, where evaluating the data for only full features wereas, limited evidence does not assumed these to does room of the data for only full features.	Since	
Open en Antonanpie Repair*	Being relation apport to difference is long term (*1 year) paint reported externet or cell basiss and between you and attractive traces between attractive cells to being	Georg	

A W35 Roommadulina	\$wy	Chatteri-e	Contai	361030-78	Chade
Anne Ogligeio					
Adaptation of Souli to Median Fall-Thickness Team	An bolt physical through and operative transmit to commanded for small to mediant fold-thickness tensor coll team?	с.	67	0.	e.
Alanapament of High Chails Partiel Trickness Learn	Is any union to a full declarers new recommended for particula- with Nath-aready method symper coll Supp."	E.+	0.		
Long-turn Non-operative Management	Is long-turn measurable management momented for priority for long-full thickness enters cull turn?	0	0	1-	p-
Puppois (Cloud	Is a clinical exem nonmendeal to dispose patients with rotator coll heard.	0	0	c-	0-
Nupoois (Imaging)	Is imaging, such as MRI, MRIA, and ultrasound, to empended to design and an MRIA and an employed "	01	0	e-	6-
Post-operative Multilization Tanàn	Advisely and delegal technication recommended for anyone what advances of senal?"	01	En-	P-	¥-
Prognostic Fastors (Apr)	Is shirt up associated with higher to new ones after sensor call	E.+	с.	с.	0
Property Factors (Walker's	Is under to compare also associated with higher to have sales	Di-	22		12
Bology Aggregation with Patrice Berlinal Products	b. Indepied supmariation of rotator call space with planker derived produce, successful for improving pairon sported automore?	1.	I+	14	1.
Xingle Barn vs. Daable Barn Repair	An hole for single row and deality row report techniques, measure-adeal for improving partient reported sourcemes after source coff result*	C>	0	0	0
Kingle-Bara vs. Deable-Bara Barairo Sat Bo-Yana	An bolk for single one and deality one report moleculars, programming for proceeding to many after strategy off people"	P.	£74		1.
Open vs. Anthensorpis Kapate	An little for open and a theorem, traine call reput techniques recommended for improving param reported subcases, range of motion or mate?	0	0	0	0
Operative Management	Is questive management of triater coll team recommended to improve concerns compared to physical theory and athenial strater coll require?	P.	0	P	- 10
Accomigizity and Resear Coll Repub	Is the matter use of accomplicity during arbitraryin states cull typic recommended for patients with small to mediant shed attend off mat?	P	E*	£1	0