

Comparing AAOS Appropriate Use Criteria with ChatGPT-4.5 and Gemini 2.0 Pro Recommendations for Hip Osteoarthritis

Charlene W Cai, Alexander Yu, Aneesh P Reddy, Jonathan J Huang, Prabhjot Singh, Samuel Kang-Wook Cho

INTRODUCTION: Hip osteoarthritis presents a complex clinical challenge, with treatment decisions influenced by a multitude of patient-specific factors, such as age, pain severity, radiographic findings, range of motion, and risk of adverse outcomes. The American Academy of Orthopaedic Surgeons (AAOS) developed appropriate use criteria (AUC) to guide treatment decisions for hip osteoarthritis based on expert consensus. This study evaluates the ability of two advanced AI (Artificial Intelligence) models, Chat Generative Pre-trained Transformer-4.5 (ChatGPT-4.5) and Gemini 2.0 Pro, to generate treatment recommendations for hip osteoarthritis by comparing their appropriateness scores for hip osteoarthritis treatment with those from the AUC.

METHODS: The AAOS AUC, ChatGPT-4.5, and Gemini 2.0 Pro were utilized to assess treatment suggestions for hip osteoarthritis. The 270 patient scenarios were assessed using a 1 to 9 scale to rate the appropriateness of various treatment options, including conservative options (risk factor assessment, activity modification, assistive devices, oral medications, intraarticular steroids, and physical therapy) and surgical interventions (arthroplasty, hip preservation, and arthrodesis). ChatGPT-4.5 and Gemini 2.0 Pro were queried to rate the treatments on the identical scale. Mean absolute and squared errors were calculated to assess the models' predictive accuracy compared to AAOS guidelines. Spearman correlations and paired t-tests ($\alpha < .05$) were conducted to evaluate consensus, with heatmaps to visualize findings.

RESULTS: Both ChatGPT-4.5 and Gemini 2.0 Pro provided hip osteoarthritis treatment recommendations that differed significantly from the AAOS guidelines ($p < 0.001$) (Tables 1 and 3). For most conservative treatments, Gemini demonstrated smaller mean absolute errors (e.g., assistive devices: +1.10 vs +2.30) compared to ChatGPT-4.5 (Tables 1 and 3). Relative to Gemini, ChatGPT-4.5 predictions showed greater standard deviations and variability in response to nuanced patient-specific factors (Tables 1 and 3; Figures 1 and 2). Though, this increased sensitivity to individual traits may also introduce potential inconsistencies in clinical decision-making. Conversely, Gemini's more conservative bias, reflected by its tighter clustering (Figures 1 and 2), may demonstrate limited responsiveness to patient-specific variation and inability to differentiate complex patient scenarios compared to ChatGPT-4.5 (Tables 1 and 3). Both models consistently underestimated the appropriateness of surgical interventions. However, the AI models demonstrated stronger rank correlations with AAOS guidelines for surgical interventions (e.g., GPT-4.5: $\rho = 0.648$ for arthroplasty; Gemini: $\rho = 0.522$) than for conservative treatments (e.g., GPT-4.5: $\rho = -0.00813$ for risk factor assessment; Gemini: $\rho = -0.129$) (Tables 2 and 4).

DISCUSSION AND CONCLUSION: Both ChatGPT-4.5 and Gemini 2.0 Pro significantly differed from AAOS recommendations in treating hip osteoarthritis ($p < 0.001$ across all treatments). While there was moderate concordance in certain treatments, the general trend towards favoring more conservative approaches and underrecommending surgical interventions raises concern about the reliability of AI-generated recommendations for clinical decision-making. If employed without human oversight, the AI models' underestimation of surgical appropriateness could lead to the undertreatment of hip osteoarthritis. Additionally, consistently over-/underestimating conservative treatments may lead to suboptimal early management decisions. As such, future development of AI models for orthopedic decision support should incorporate specialized training on clinical guidelines and real-world patient outcomes to improve both accuracy and applicability in practice. Currently, these models are not suited to generate recommendations for patients independently, but rather assist decision-making based on validated training data.

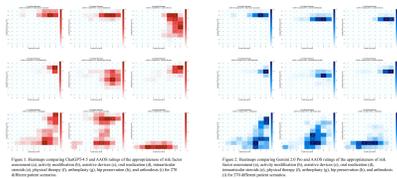


Table 1: Mean absolute errors (MAE) and Spearman's rho (ρ) for ChatGPT-4.5 and Gemini 2.0 Pro across 14 treatments for 270 patient scenarios.

Treatment	Mean Error	Mean Absolute Error	Mean Squared Error	ρ	P-value
Risk Factor Assessment	-0.99(1.12)	1.00(1.12)	1.00(1.12)	-0.00813	<.001
Activity Modification	-1.00(1.12)	1.00(1.12)	1.00(1.12)	-0.00813	<.001
Assistive Devices	-1.10(1.12)	1.10(1.12)	1.10(1.12)	-0.00813	<.001
Oral Medications	-1.10(1.12)	1.10(1.12)	1.10(1.12)	-0.00813	<.001
Intraarticular Steroids	-1.10(1.12)	1.10(1.12)	1.10(1.12)	-0.00813	<.001
Physical Therapy	-1.10(1.12)	1.10(1.12)	1.10(1.12)	-0.00813	<.001
Arthroplasty	-1.10(1.12)	1.10(1.12)	1.10(1.12)	-0.00813	<.001
Hip Preservation	-1.10(1.12)	1.10(1.12)	1.10(1.12)	-0.00813	<.001
Arthrodesis	-1.10(1.12)	1.10(1.12)	1.10(1.12)	-0.00813	<.001

Table 2: Spearman's rho (ρ) for ChatGPT-4.5 and Gemini 2.0 Pro across 14 treatments for 270 patient scenarios.

Treatment	Spearman's rho	P-value
Risk Factor Assessment	-0.00813	<.001
Activity Modification	-0.00813	<.001
Assistive Devices	-0.00813	<.001
Oral Medications	-0.00813	<.001
Intraarticular Steroids	-0.00813	<.001
Physical Therapy	-0.00813	<.001
Arthroplasty	0.648	<.001
Hip Preservation	0.648	<.001
Arthrodesis	0.648	<.001

Table 3: Spearman's rho (ρ) for Gemini 2.0 Pro across 14 treatments for 270 patient scenarios.

Treatment	Mean Error	Mean Absolute Error	Mean Squared Error	ρ	P-value
Risk Factor Assessment	-0.99(1.12)	1.00(1.12)	1.00(1.12)	-0.129	<.001
Activity Modification	-1.00(1.12)	1.00(1.12)	1.00(1.12)	-0.129	<.001
Assistive Devices	-1.10(1.12)	1.10(1.12)	1.10(1.12)	-0.129	<.001
Oral Medications	-1.10(1.12)	1.10(1.12)	1.10(1.12)	-0.129	<.001
Intraarticular Steroids	-1.10(1.12)	1.10(1.12)	1.10(1.12)	-0.129	<.001
Physical Therapy	-1.10(1.12)	1.10(1.12)	1.10(1.12)	-0.129	<.001
Arthroplasty	-1.10(1.12)	1.10(1.12)	1.10(1.12)	0.522	<.001
Hip Preservation	-1.10(1.12)	1.10(1.12)	1.10(1.12)	0.522	<.001
Arthrodesis	-1.10(1.12)	1.10(1.12)	1.10(1.12)	0.522	<.001

Table 4: Spearman's rho (ρ) for Gemini 2.0 Pro across 14 treatments for 270 patient scenarios.

Treatment	Spearman's rho	P-value
Risk Factor Assessment	-0.129	<.001
Activity Modification	-0.129	<.001
Assistive Devices	-0.129	<.001
Oral Medications	-0.129	<.001
Intraarticular Steroids	-0.129	<.001
Physical Therapy	-0.129	<.001
Arthroplasty	0.522	<.001
Hip Preservation	0.522	<.001
Arthrodesis	0.522	<.001

Figure 1: Heatmaps comparing ChatGPT-4.5 and Gemini 2.0 Pro ratings of the appropriateness of 14 treatments across 270 patient scenarios. The heatmaps show a distribution of scores from 1 (blue) to 9 (red).

Figure 2: Heatmaps comparing Gemini 2.0 Pro and AAOS ratings of the appropriateness of 14 treatments across 270 patient scenarios. The heatmaps show a distribution of scores from 1 (blue) to 9 (red).