

# A Comparison Between the Efficacy of ChatGPT, Grok, and Claude Sonnet in Analyzing Common Sports-Related Radiologic Imaging

Michael Lane Moore, Ahmad R Alhankawi, Collin Braithwaite, Alex Miguel Holle, Anikar Chhabra

**INTRODUCTION:** The use of artificial intelligence (AI) has been a subject of major interest in the field of medicine. Previous studies have shown AI's powerful capabilities in the accuracy of medical inquiries, but a lack of data is published using radiologic images. The purpose of this paper is two-fold: (1) to assess the general ability for AI to diagnose common sport-related pathologies using radiologic imaging, and (2) to compare ChatGPT with two of its competitors, Grok and Claude Sonnet.

**METHODS:** ChatGPT 4.0, Grok 2, and Claude 3.5 Sonnet were utilized. Five common orthopedic sports pathologies were chosen: anterior cruciate ligament (ACL) tears, posterior cruciate ligament (PCL) tears, meniscal tears, chondral pathologies, and rotator cuff tears. Fifty images representing each pathology were randomly collected from a radiologic imaging database, when possible, which included radiographic images, computed tomography (CT), and magnetic resonance imaging (MRI). Normal images were collected that corresponded to each diagnostic category. A new query was used for each image to assess a likely diagnosis. Receiver operator characteristic curves and area under the curve values were calculated to assess the accuracy of each AI platform.

**RESULTS:** Regarding diseased images, ChatGPT, Grok and Claude Sonnet accurately identified the pathology in 23%, 16%, and 17% of images, respectively. ChatGPT and Grok were most accurate at identifying meniscus pathologies (ChatGPT: 48%, Grok: 42%), while Claude Sonnet was most accurate at identifying ACL pathologies (31%). The AUC for ChatGPT, Grok, and Claude Sonnet was 0.21, 0.16, and 0.16, respectively. There were no differences in performance between the three platforms overall or within any of the diagnostic categories. ChatGPT and Claude Sonnet had a statistically higher rate of correctly diagnosing diseased MRI images (ChatGPT  $p < 0.001$ , Claude 3.5 Sonnet  $p = 0.003$ ) while Grok showed no differences in accuracy among imaging modalities ( $p = 0.088$ ).

**DISCUSSION AND CONCLUSION:** Open-Sourced AI platforms ChatGPT, Grok and Claude Sonnet had a diagnostic accuracy of less than 25% for common orthopedic sports pathologies, and had AUC values well under 0.5, indicating very low accuracy. Although ChatGPT generally operated better than Grok and Claude Sonnet, the overall results were unremarkable. We do not recommend its current use in image-based diagnoses.

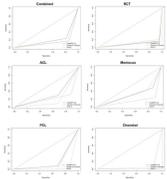


Figure 3. ROC curves generated for the overall combined results of images based on ChatGPT 4.0, Grok 2, and Claude 3.5 Sonnet across the categories of diagnostic categories.

Pathology	Total	ACL	PCL	Meniscus	Chondral	Rotator Cuff
ChatGPT	50	14	4	14	6	5
Grok	23	20	0	10	0	10
Claude	31	1	0	13	1	8
AUC	0.21	0.2	0.2	0.2	0.2	0.2

Pathology	ChatGPT	Grok	Claude
ACL	0.23	0.16	0.17
PCL	0.21	0.16	0.16
Meniscus	0.48	0.42	0.17
Chondral	0.21	0.16	0.16
Rotator Cuff	0.21	0.16	0.16

Pathology	ChatGPT	Grok	Claude
ACL	0.23	0.16	0.17
PCL	0.21	0.16	0.16
Meniscus	0.48	0.42	0.17
Chondral	0.21	0.16	0.16
Rotator Cuff	0.21	0.16	0.16

Pathology	ChatGPT	Grok	Claude
ACL	0.23	0.16	0.17
PCL	0.21	0.16	0.16
Meniscus	0.48	0.42	0.17
Chondral	0.21	0.16	0.16
Rotator Cuff	0.21	0.16	0.16

Pathology	ChatGPT	Grok	Claude
ACL	0.23	0.16	0.17
PCL	0.21	0.16	0.16
Meniscus	0.48	0.42	0.17
Chondral	0.21	0.16	0.16
Rotator Cuff	0.21	0.16	0.16