

Evaluation of ChatGPT in Assessing Postoperative Outcomes Following Flexor Tendon Repair (FTR)

Kismat A Touhid, Firdavs Kurbanov, Amy Phan, Constantinos Ketonis

INTRODUCTION: Large language models (LLMs) like ChatGPT have the potential to ameliorate the administrative burden hand surgeons face when querying the electronic medical record (EMR) for patient data. Our pilot case series evaluated the reliability of ChatGPT-4 in synthesizing rehabilitation data from postoperative EMR notes following FTR to assess its utility and limitations.

METHODS: Postoperative progress notes for 5 FTR patients were anonymized and input into ChatGPT-4 Turbo with memory disabled. Using a standardized multi-prompt framework, we directed the model to: (1) generate qualitative summary of patient recovery, (2) extract and tabulate quantitative metrics from postoperative notes, and (3) graph the data for easy visualization. The output was then cross-referenced against original notes within the EMR for accuracy and omissions.

RESULTS: Analysis of 95 clinical notes (29,000 words) across postoperative office, physical therapy (PT), and occupational therapy (OT) documentation revealed that GPT-4 achieved sufficient coverage of postoperative data without generating hallucinated claims. However, qualitative summaries demonstrated contextual omissions (e.g., social determinants, injury etiology) in 4/5 patients (Table 1, Figure 1). Quantitative metrics, including joint-specific ranges of motion (ROM), grip strength, and pain scores, were tabulated accurately and with precision; however, graphical outputs exhibited inconsistencies: omitted extension ROM trajectories in 2/5 patients, omitted total active motion (TAM) in 1/5 patients, omitted non-injured limb grip strength in 2/5 patients, and misrepresentation of patient-reported pain ranges in 2/5 patients (Table 2).

DISCUSSION AND CONCLUSION:

Integrating structured prompting frameworks allows ChatGPT-4 to reliably aggregate postoperative rehabilitation data with high fidelity. However, its limitations in contextual granularity and graphical representation necessitate clinician-guided oversight to mitigate the risk of oversimplification in complex recovery trajectories.

Accuracy Assessment of GPT-4 Qualitative Outputs Across 5 Patients

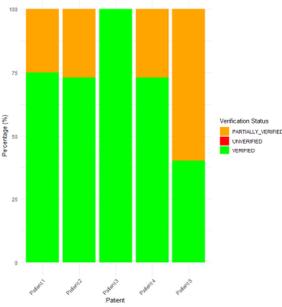


Table 1. Accuracy Assessment of GPT-4 Outputs: Verification Rates for Qualitative Data Across Five Flexor Tendon Repair Patients

Patient ID (Clinical Reference)	Total # of Qualitative Claims based on Clinical Record Prompting	% of Qualitative Claims Verified against Clinical Record (VERIFIED)	% of Qualitative Claims with Potential Contextual Omissions (PARTIALLY VERIFIED)	% of Qualitative Claims Unverified against Clinical Record (NOT VERIFIED)
Patient 1 (Table A1)	12	75% (9/12)	24% (3/12)	0% (0/12)
Patient 2 (Table A2)	15	73% (11/15)	27% (4/15)	0% (0/15)
Patient 3 (Table A3)	12	100% (12/12)	0% (0/12)	0% (0/12)
Patient 4 (Table A4)	10	40% (4/10)	60% (6/10)	0% (0/10)

Table 2. Accuracy Assessment of GPT-4 Outputs: Verification Rates for Quantitative Data Across Five Flexor Tendon Repair Patients

Patient ID	Accuracy of Data Metrics	Graphical Issues (Omissions/Hallucinations)
Patient 1 (Table A1, Figure 1)	MPF Extension/Elevation: 10/10 PPF Extension/Flexion: 10/10 Grip Strength: 10/10 Pain Score: 10/10 ROM: 10/10 OT Non-Operated Limb: 10/10	MPF Extension Graph omitted. Grip Strength Value Graph for Non-Operated Limb omitted.
Patient 2 (Table A2, Figure A2)	MPF Extension/Elevation: 10/10 PPF Extension/Flexion: 10/10 Grip Strength: 10/10 Pain Score: 10/10 ROM: 10/10 OT Non-Operated Limb: 10/10	Pain Score Ranges omitted (instead of reporting Full Range).
Patient 3 (Table A3, Figure A3)	MPF Extension/Elevation: 10/10 PPF Extension/Flexion: 10/10 Grip Strength: 10/10 Pain Score: 10/10 ROM: 10/10 OT Non-Operated Limb: 10/10	MPF/PPF/OT Extension Graphs Omitted.
Patient 4 (Table A4, Figure A4)	MPF Extension/Elevation: 10/10 PPF Extension/Flexion: 10/10 Grip Strength: 10/10 Pain Score: 10/10 ROM: 10/10 OT Non-Operated Limb: 10/10	MPF/PPF/OT Extension Graphs Omitted. Pain Score Ranges (instead of only Lower Bound). Grip Strength Value Graph for Non-Operated Limb Omitted.
Patient 5 (Table A5, Figure A5)	MPF Extension/Elevation: 10/10 PPF Extension/Flexion: 10/10 Grip Strength: 10/10 Pain Score: 10/10 ROM: 10/10 OT Non-Operated Limb: 10/10	Total Active Motion (TAM) Graph omitted.