# Adapted large language models can outperform medical experts in clinical text summarization
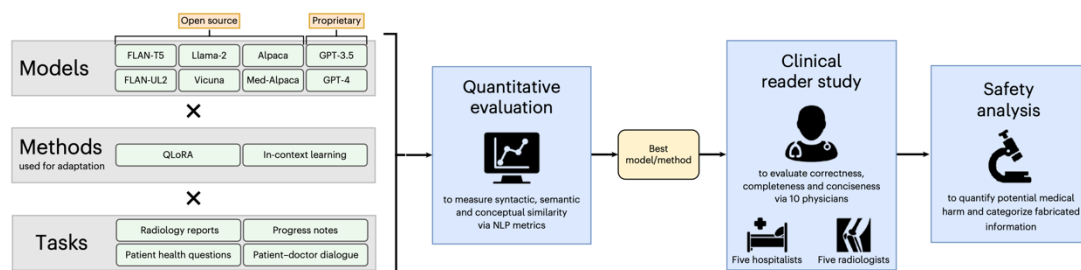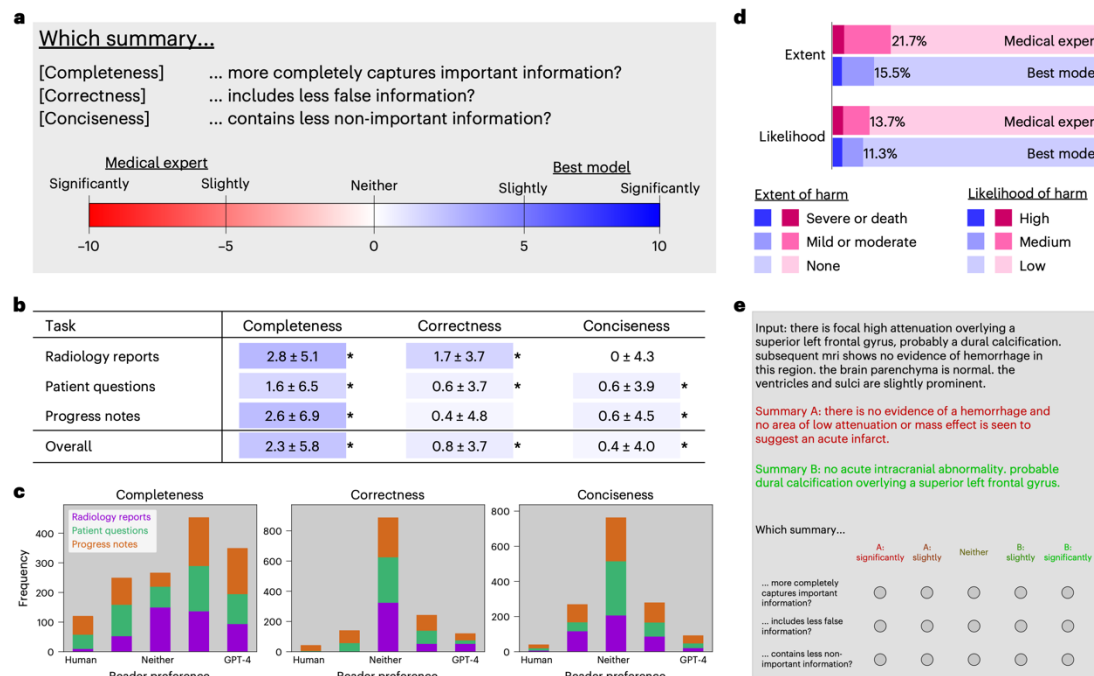
Dave Van Veen [1,2] ✉, Cara Van Uden[2,3], Louis Blankemeier[1,2], Jean-Benoit Delbrouck[2], Asad Aali[4], Christian Bluethgen [2,5], Anuj Pareek [2,6], Malgorzata Polacin[5], Eduardo Pontes Reis[2,7], Anna Seehofnerová[8,9], Nidhi Rohatgi [8,10], Poonam Hosamani[8], William Collins [8], Neera Ahuja[8], Curtis P. Langlotz [2,8,9,11], Jason Hom[8], Sergios Gatidis[2,9], John Pauly[1] & Akshay S. Chaudhari [2,9,11,12]

**Fig. 1 | Framework overview.** First, we quantitatively evaluated each valid combination (×) of LLM and adaptation method across four distinct summarization tasks comprising six datasets. We then conducted a clinical reader study in which 10 physicians compared summaries of the best model/method against those of a medical expert. Lastly, we performed a safety analysis to categorize different types of fabricated information and to identify potential medical harm that may result from choosing either the model or the medical expert summary.



**Fig. 4 | Clinical reader study. a**, Study design comparing summaries from the best model versus that of medical experts on three attributes: completeness, correctness and conciseness. **b**, Results. Highlight colors correspond to a value's location on the color spectrum. Asterisks (*) denote statistical significance by a one-sided Wilcoxon signed-rank test, *P* < 0.001. **c**, Distribution of reader scores for each summarization task across attributes. Horizontal axes denote reader preference as measured by a five-point Likert scale. Vertical axes denote frequency count, with 1,500 total cases for each plot. **d**, Extent and likelihood of possible harm caused by choosing summaries from the medical expert (pink) or best model (purple) over the other. **e**, Reader study user interface.