# RADIUS:
# Community-Driven Radiology AI Large Language Model and Vision Language Model Leaderboard

David Li, MD; Jaron Chong, MD, MHI

Schulich
MEDICINE & DENTISTRY

Western

# Disclosures

- DL: No relevant relationships.

- JC: Chair of the Canadian Association of Radiologists Artificial Intelligence Standing Committee and board member for AMS Healthcare.

# Introduction

- Exponential growth in large language models (LLMs) and vision-language models (VLMs) presents tremendous potential to transform radiology

- Yet, evaluating and comparing model performance in radiology remains challenging

# Methods

- Standardized platform for evaluating and comparing models across diverse radiological tasks and datasets, integrating both text and images

- Evaluation framework developed in collaboration with radiologists featuring domain-specific criteria

- Initialized using published results and is open to contributions from both academic and industry partners
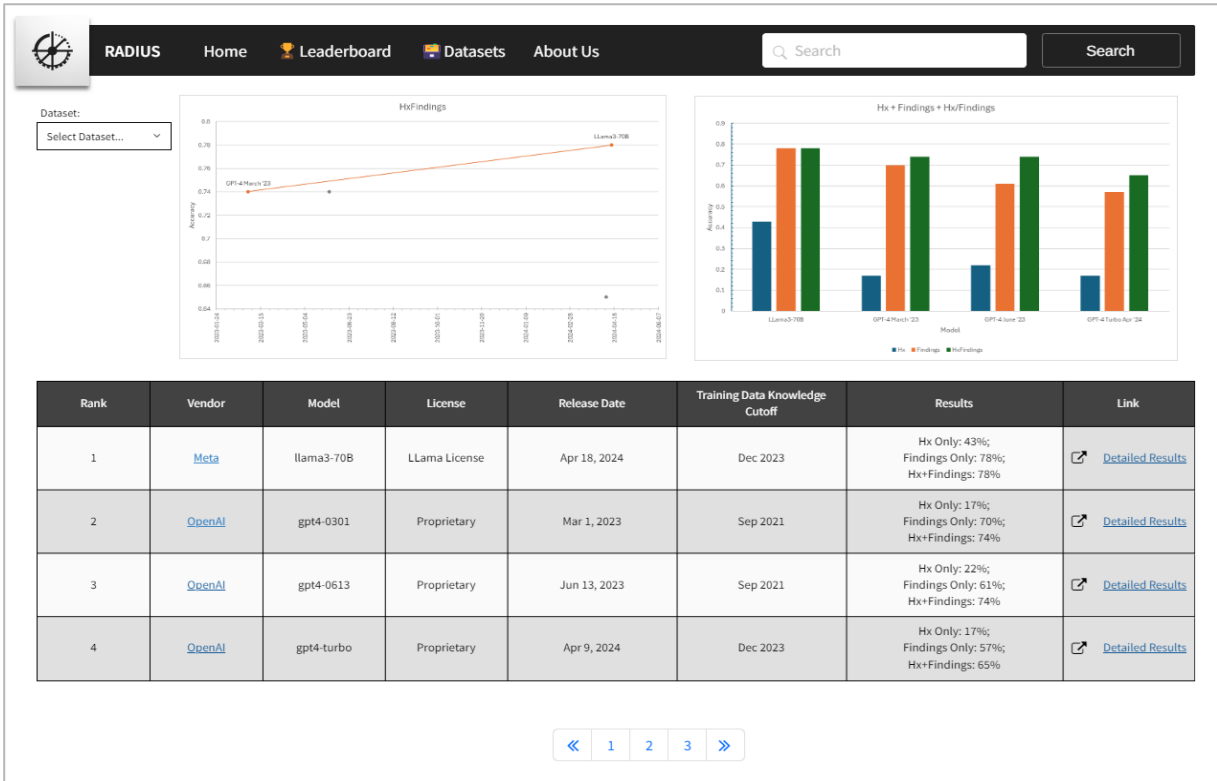
# RADIUS



Figure 1: RADIUS website for evaluation and comparison of LLMs on differential diagnoses generation task. Models are ordered by dataset-specific performance metrics plotted over time and by model.
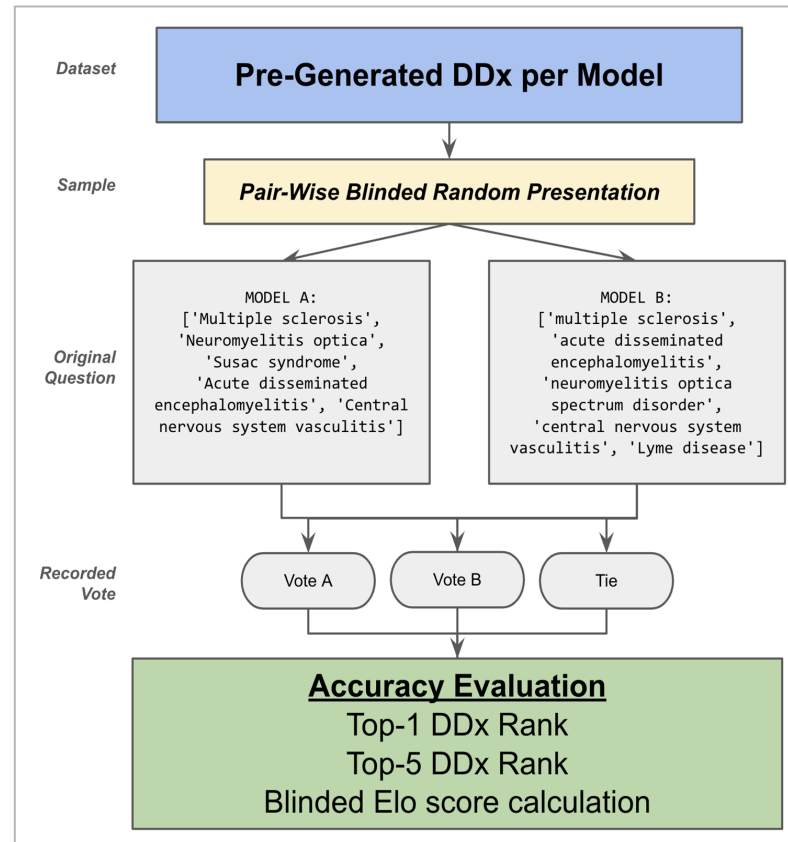
# RADIUS



Figure 2: Blinded voting system for radiologists to evaluate model performance on subjective interpretive tasks.

# Results

- Radiologists evaluated multiple proprietary and open-source models using a blinded voting system

- Model performance was transparently reported and ranked, accompanied by longitudinal analysis

- RADIUS offers a viable and effective approach for addressing the challenge of evaluating and comparing model performance in radiology

# Discussion

- RADIUS promotes fairness, transparency, and collaboration within the radiology and AI communities

- This represents a step towards standardizing model evaluation in radiology

- Ongoing development will support further quality improvement in clinical applications of generative AI

# Thank you