

## COMPREHENSIVE LLM PERFORMANCE ANALYSIS

=====

Report Generated: 2025-06-29 13:59:44

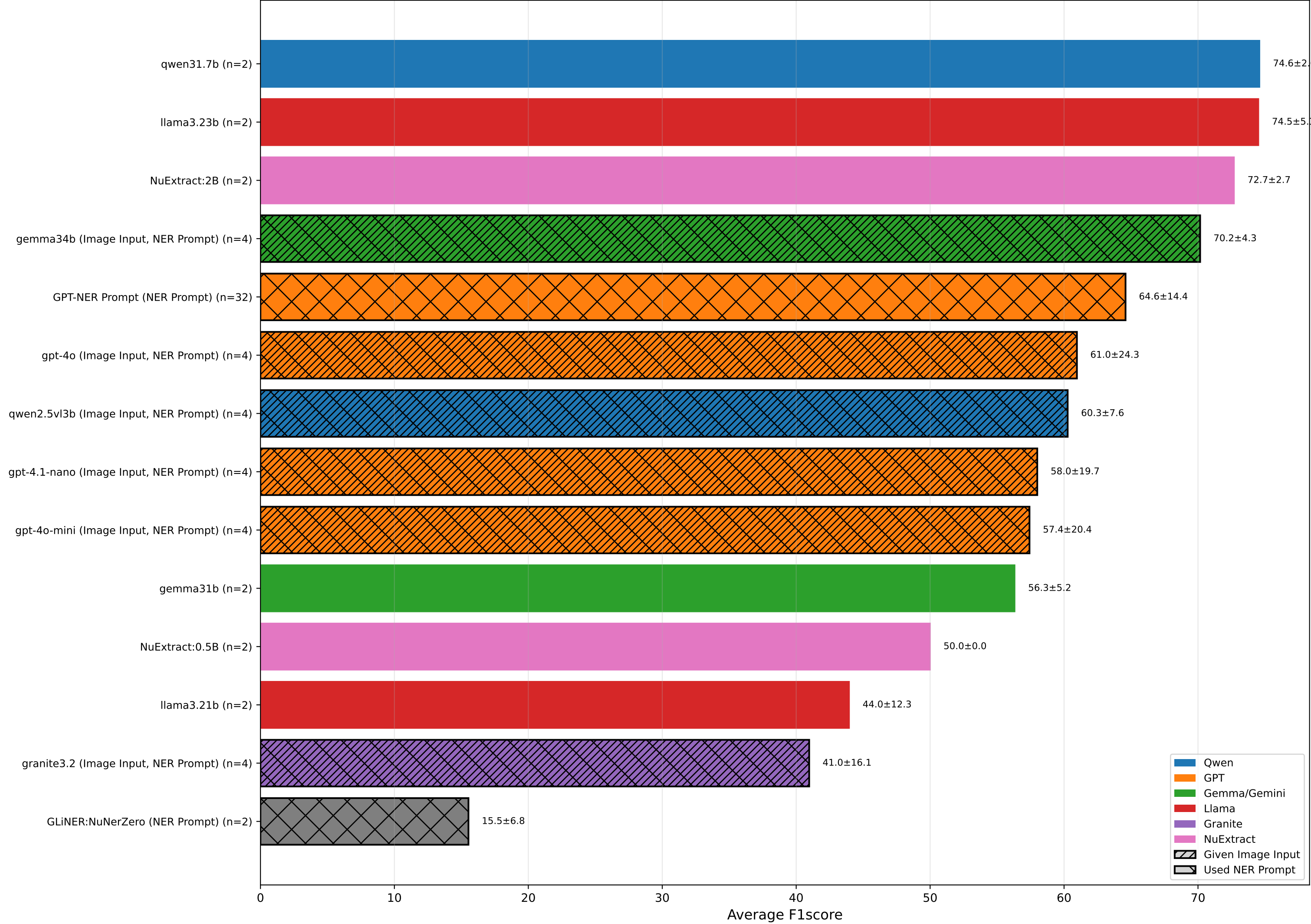
### OVERALL PERFORMANCE METRICS:

- Average F1 Score: 60.476 (Std: 17.169)
- Average Accuracy: 45.408 (Std: 16.497)
- Average Precision: 64.073 (Std: 20.554)
- Average Recall: 60.364 (Std: 17.576)

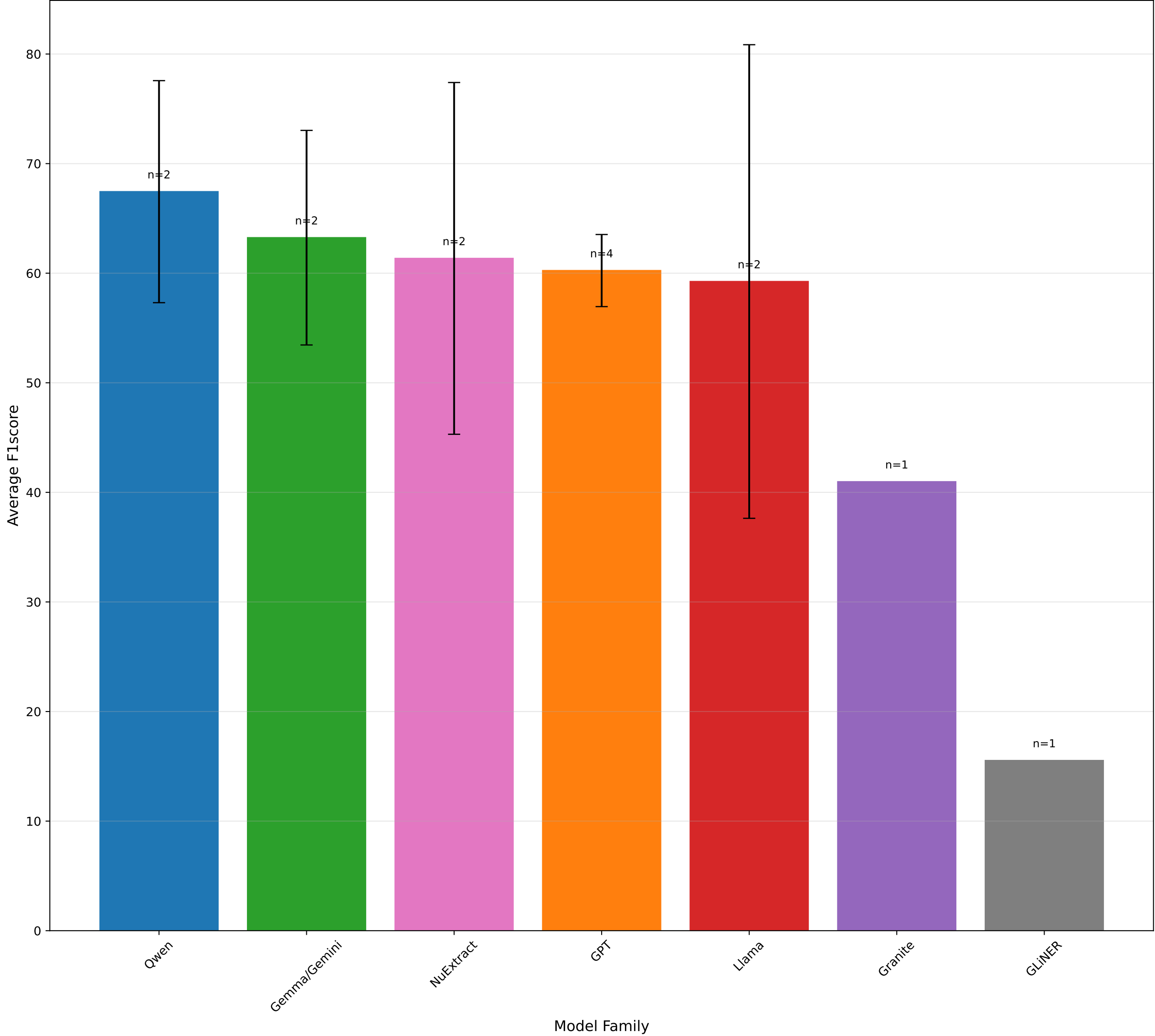
### GROUPED MODEL F1 SCORE STATISTICS:

- Unique Base Models: 14
- Total Test Instances: 70
- Best Performing Model: qwen31.7b (F1: 74.60)
- Worst Performing Model: GLiNER:NuNerZero (F1: 15.53)
- Overall Average F1: 57.14
- Models with Vision: 6

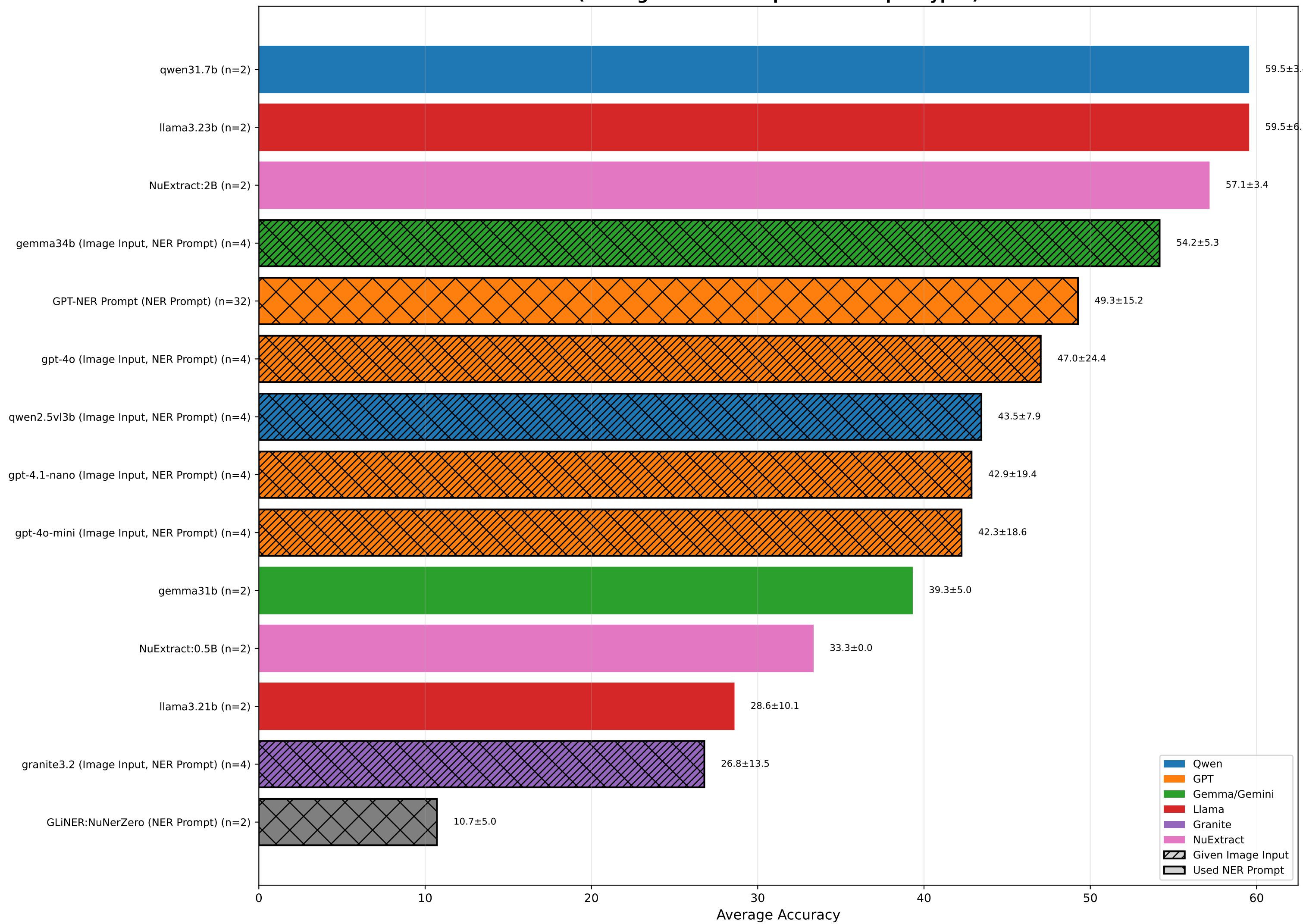
Overall F1score Performance - Models Grouped by Base Name  
(Averaged across hospitals and input types)



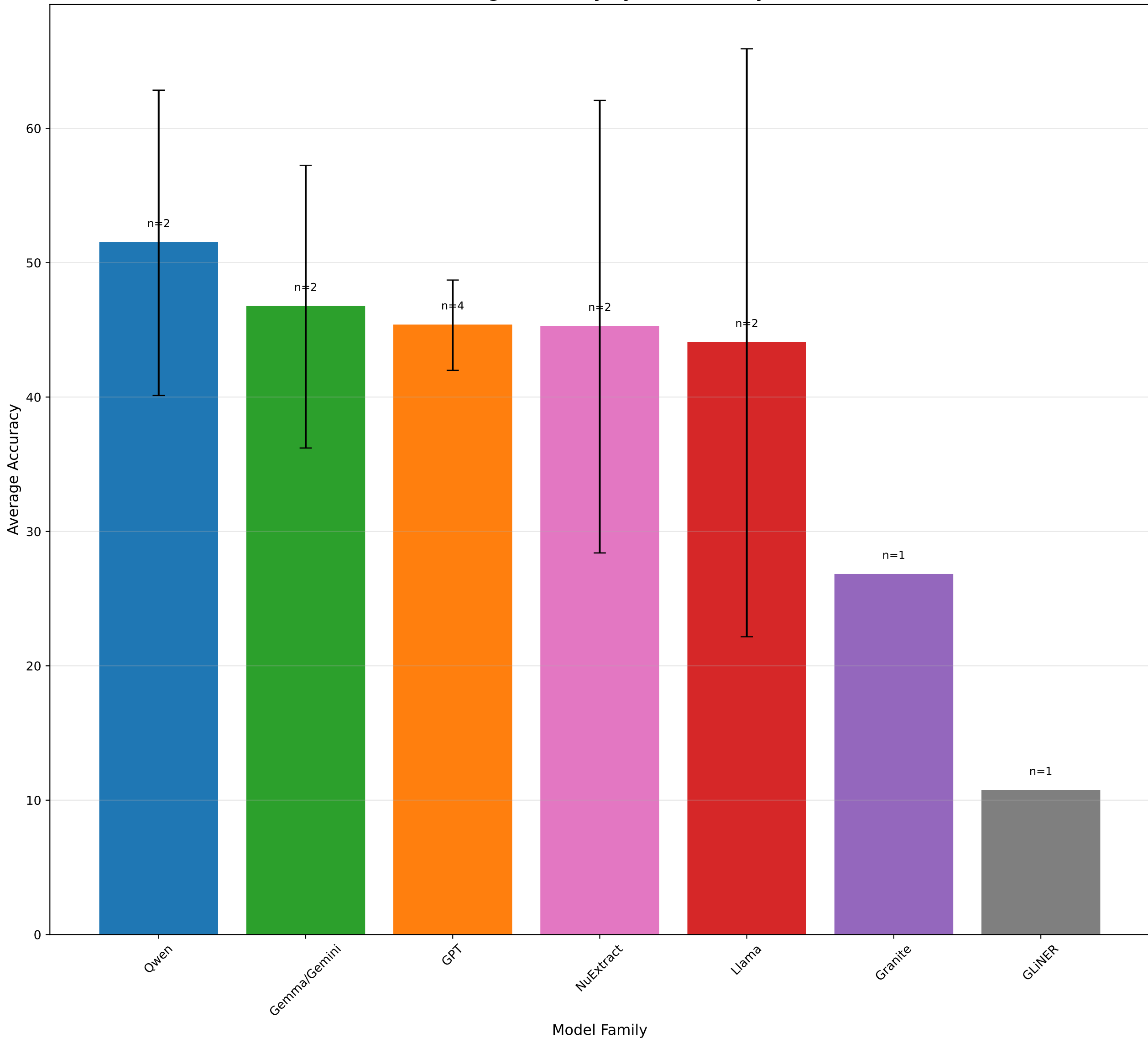
Average F1score by Model Family



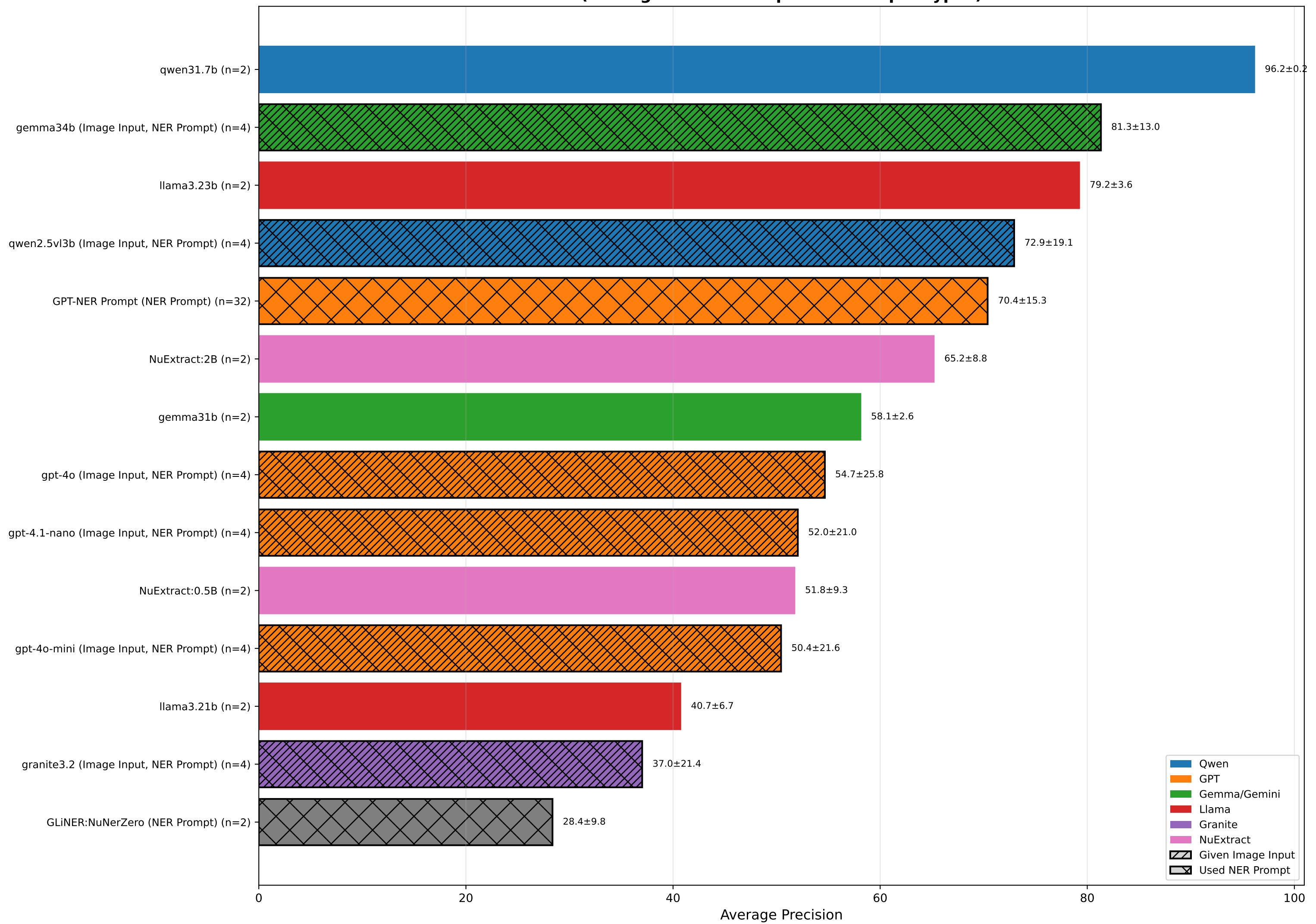
Overall Accuracy Performance - Models Grouped by Base Name  
(Averaged across hospitals and input types)



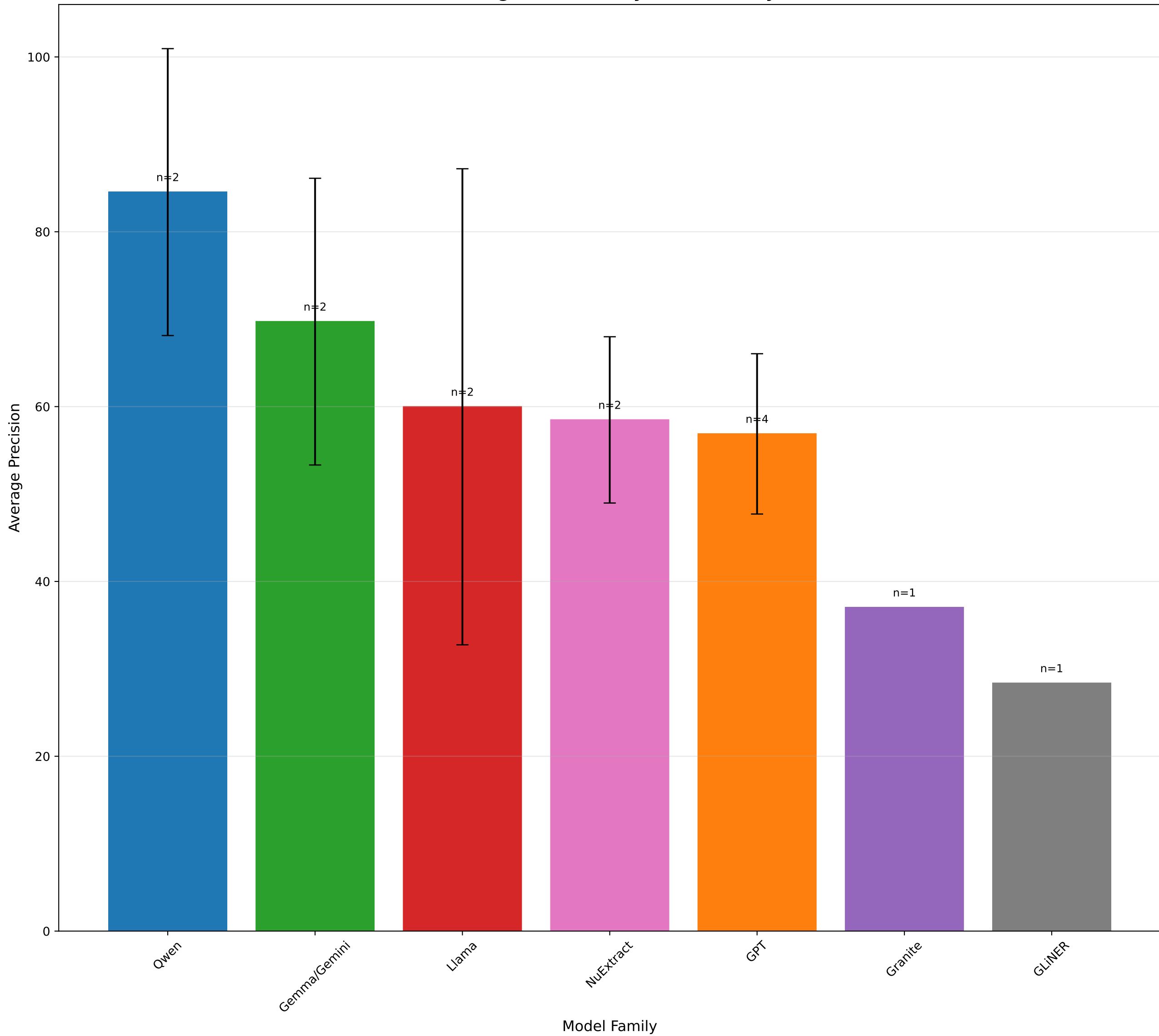
Average Accuracy by Model Family



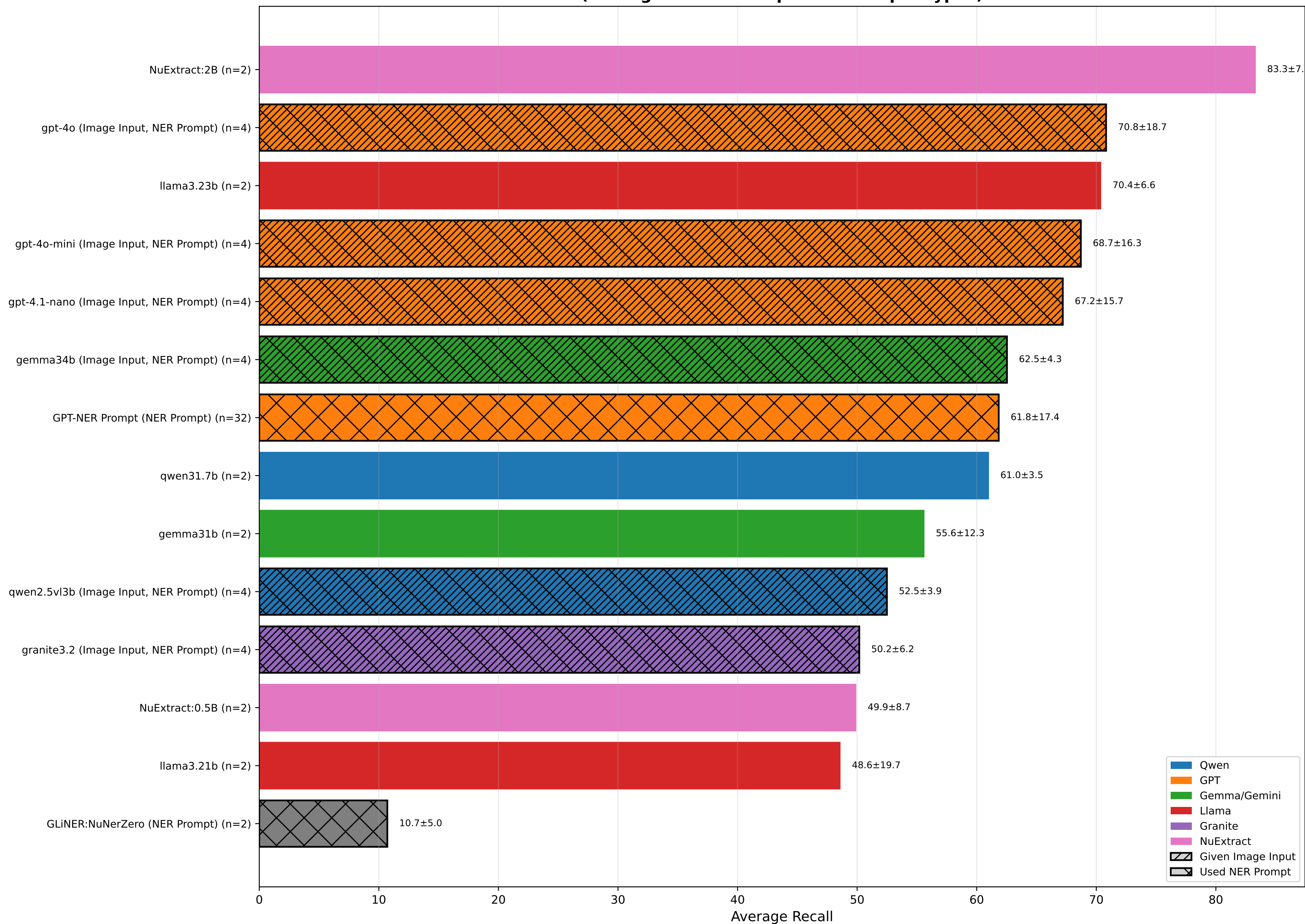
Overall Precision Performance - Models Grouped by Base Name  
(Averaged across hospitals and input types)



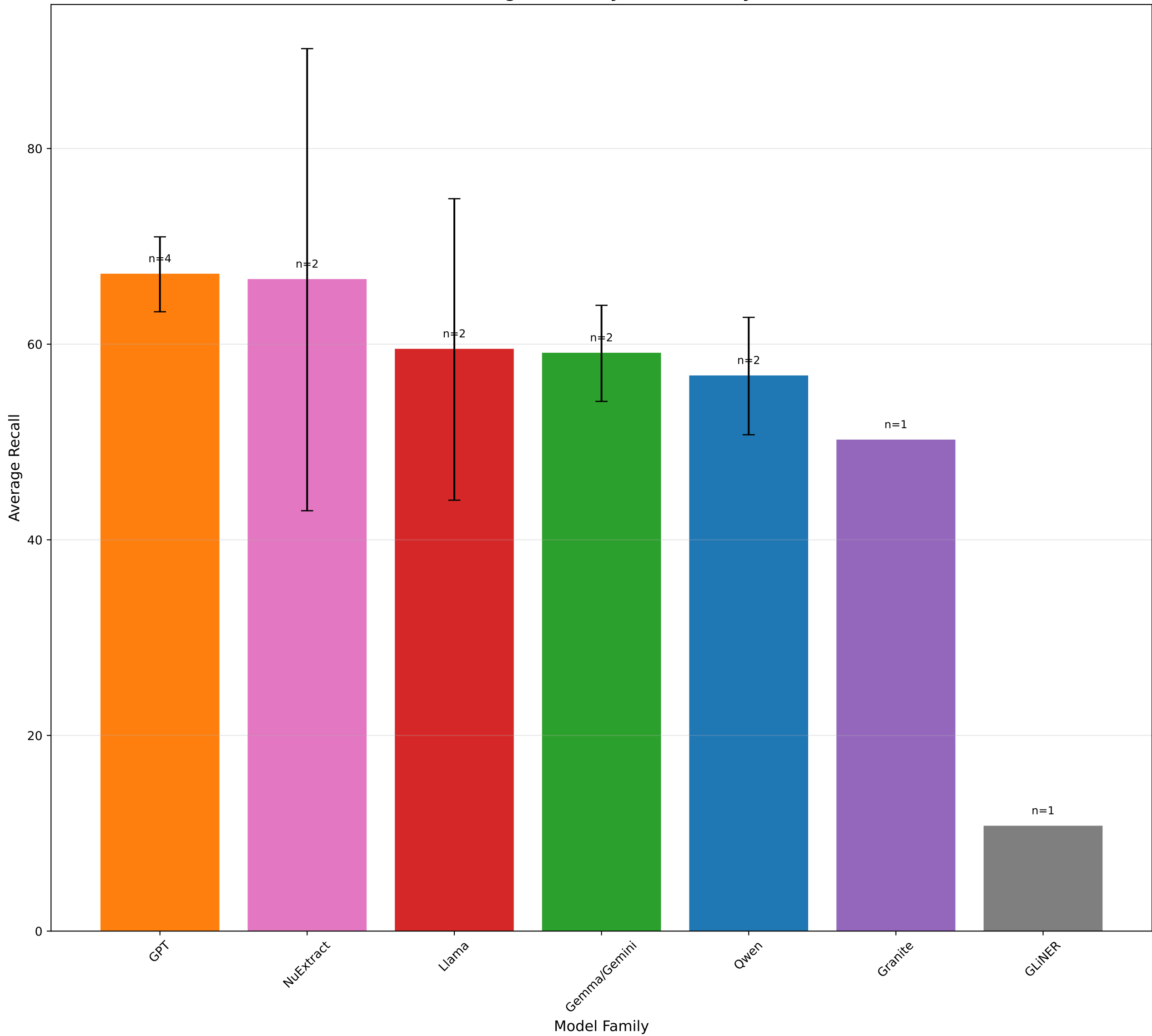
Average Precision by Model Family



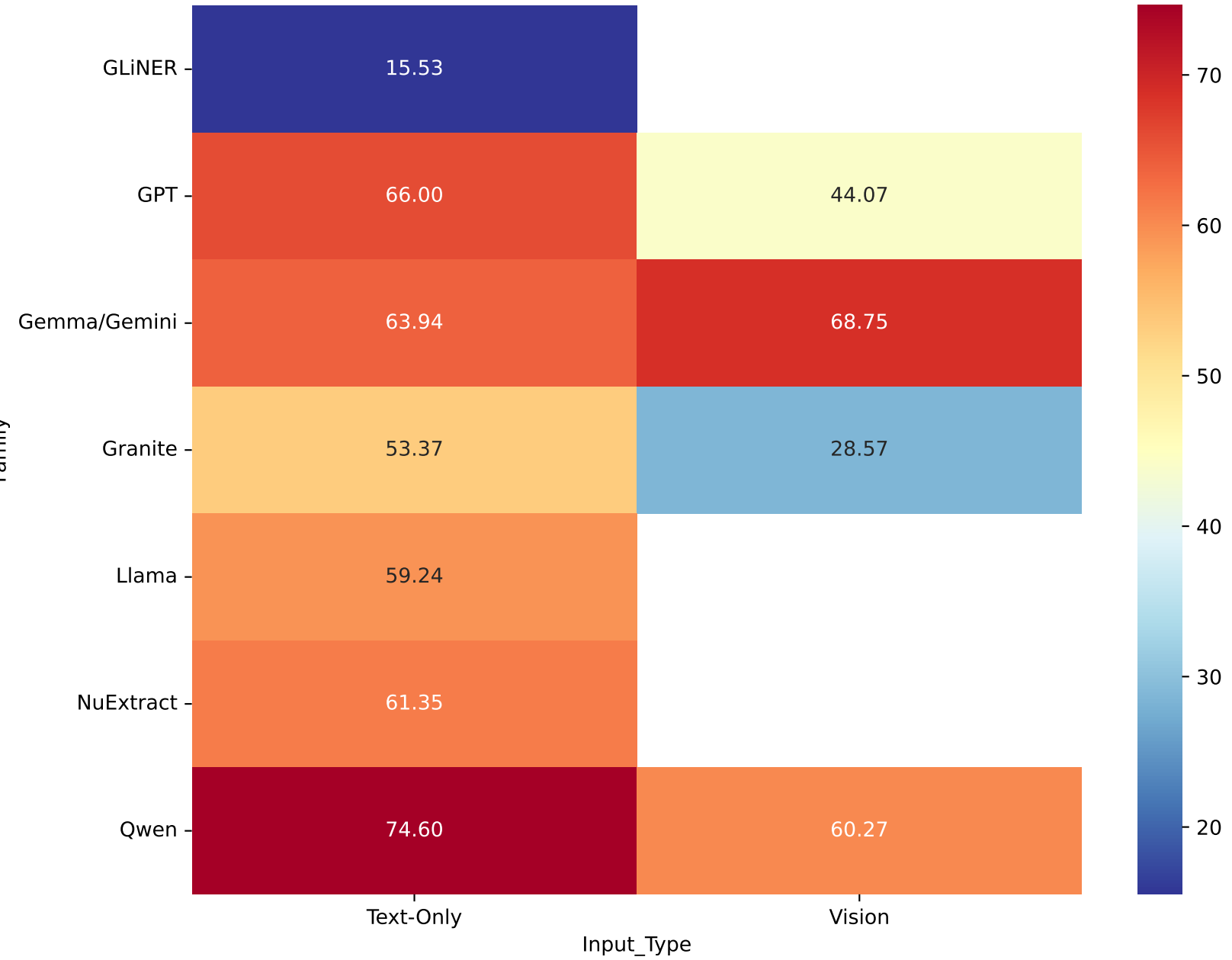
Overall Recall Performance - Models Grouped by Base Name  
(Averaged across hospitals and input types)



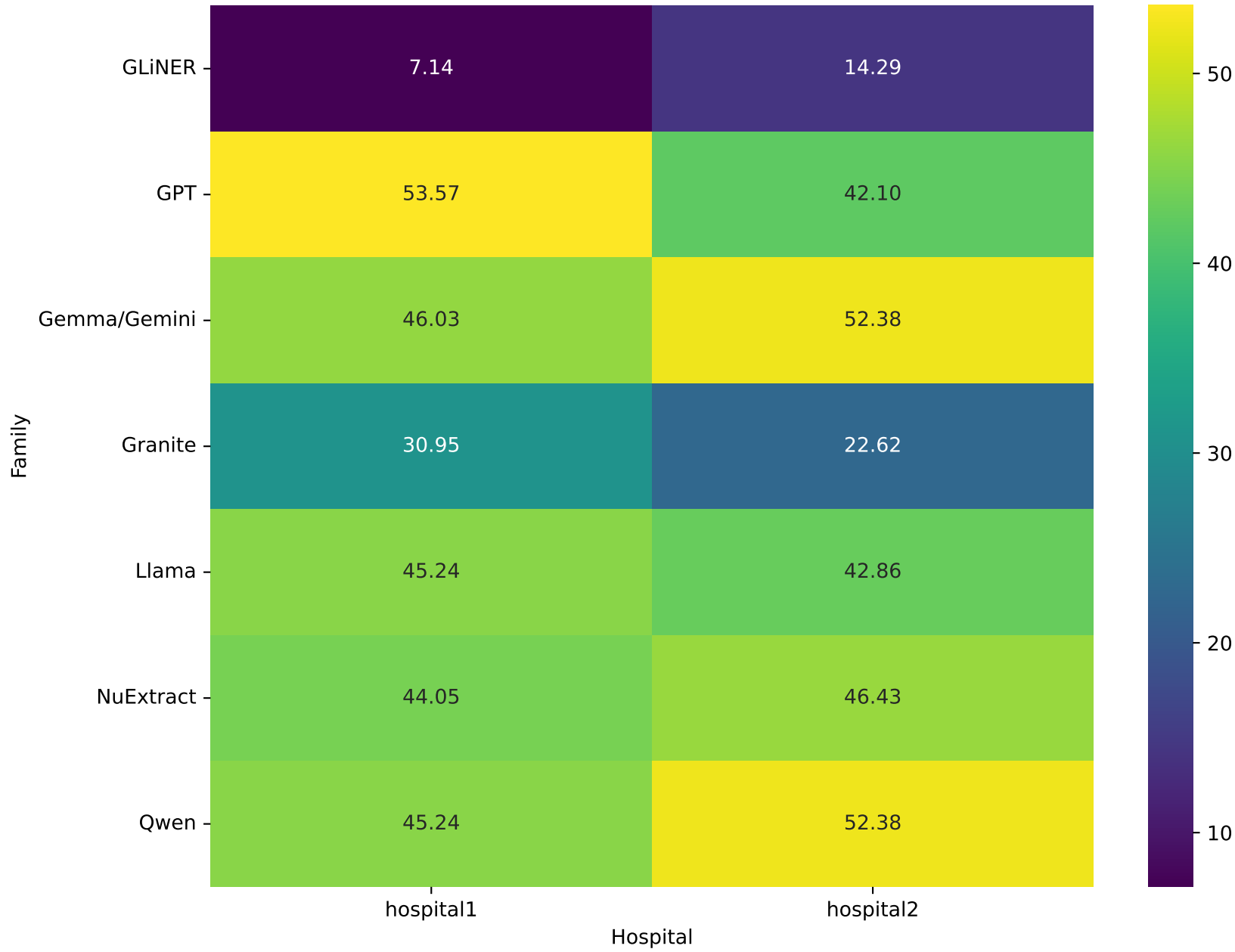
Average Recall by Model Family



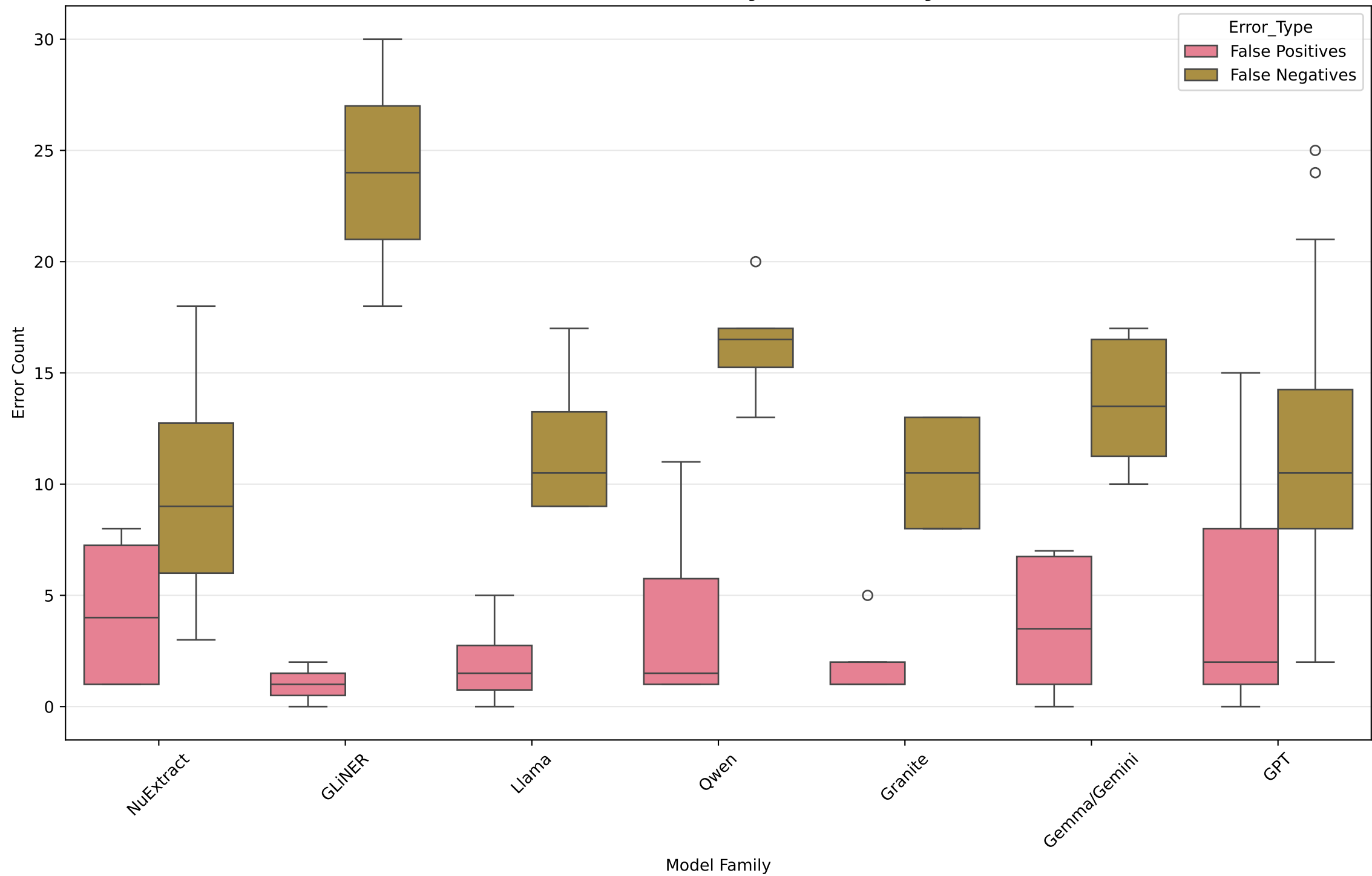
Average F1 Score: Family vs Input Type



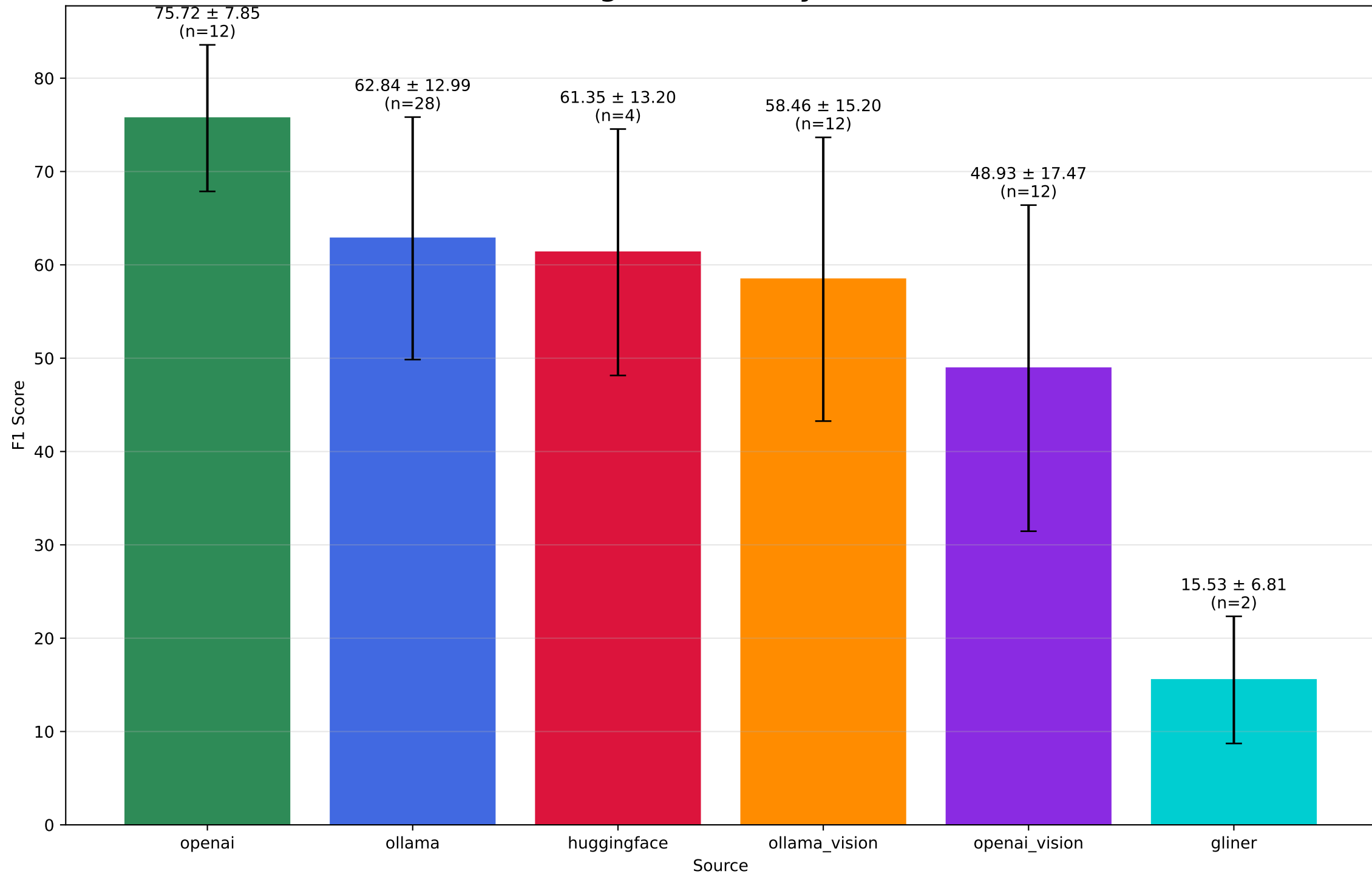
Average Accuracy: Family vs Hospital



# Error Distribution by Model Family

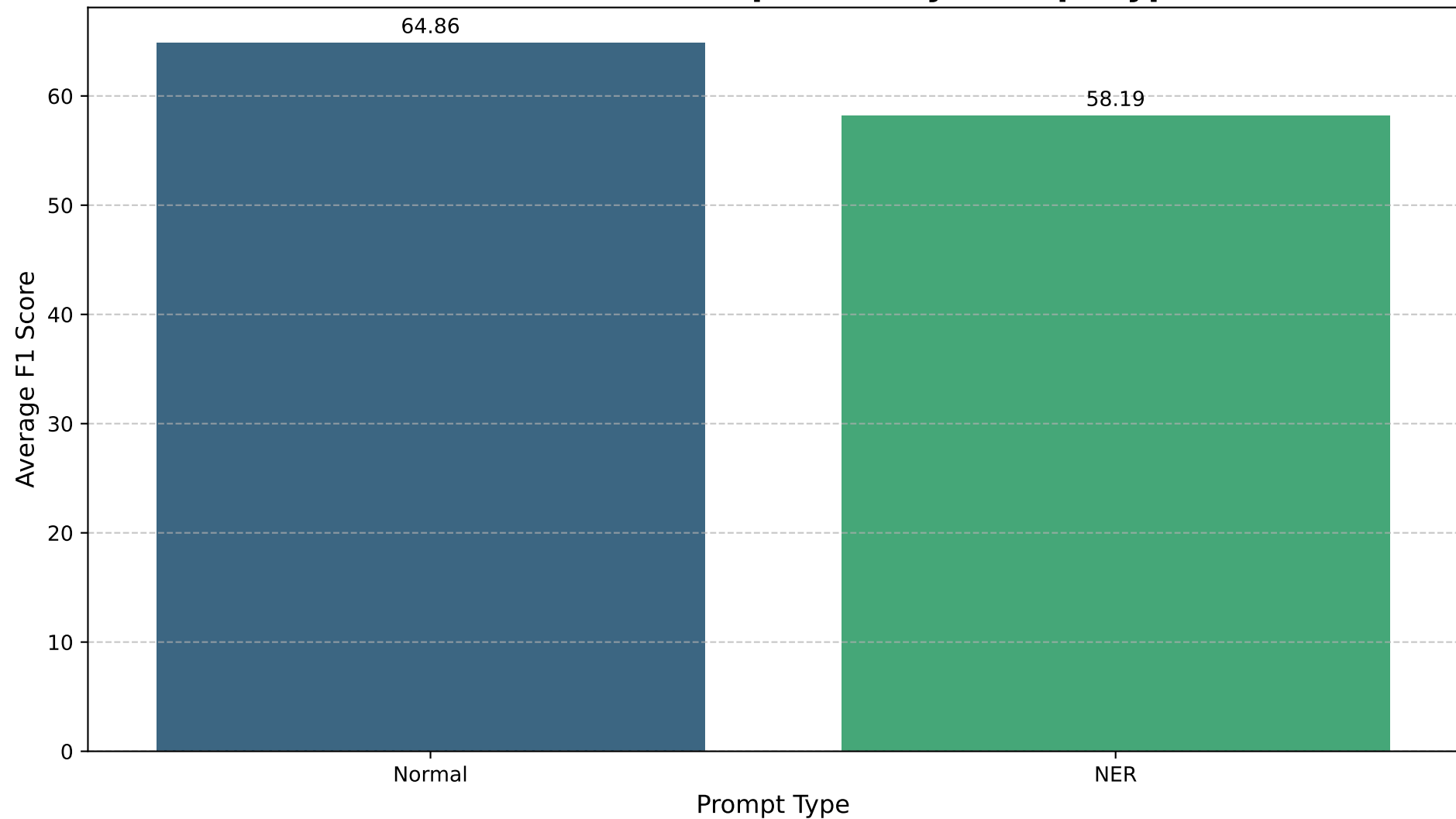


# Average F1 Score by Source

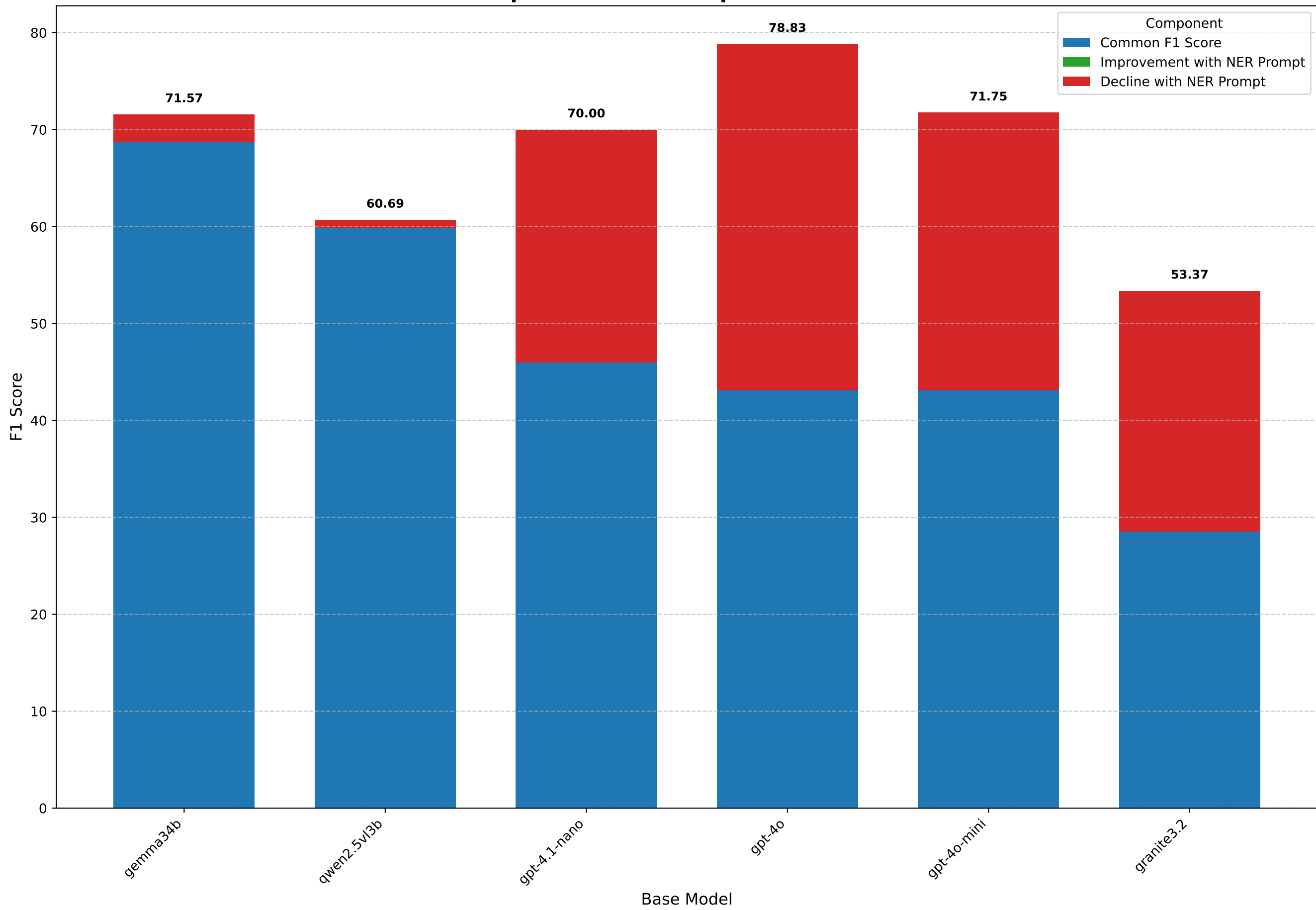




# Overall F1 Score Comparison by Prompt Type



Impact of NER Prompt on F1 Score



Impact of Image Input on F1 Score

