

## COMPREHENSIVE LLM PERFORMANCE ANALYSIS

=====

Report Generated: 2025-07-06 23:52:13

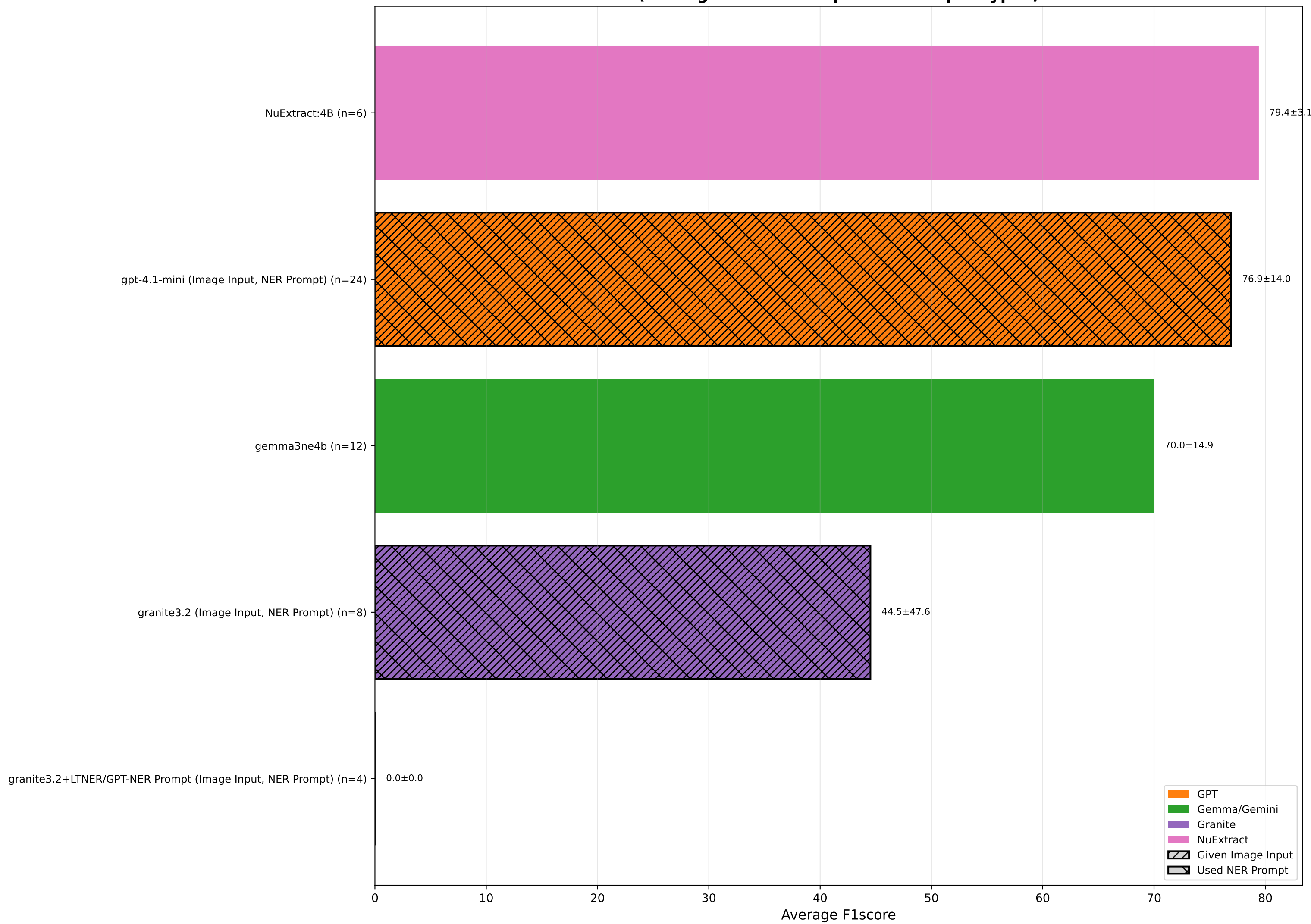
### OVERALL PERFORMANCE METRICS:

- Average F1 Score: 65.147 (Std: 30.121)
- Average Accuracy: 54.169 (Std: 27.037)
- Average Precision: 68.464 (Std: 33.019)
- Average Recall: 63.863 (Std: 30.084)

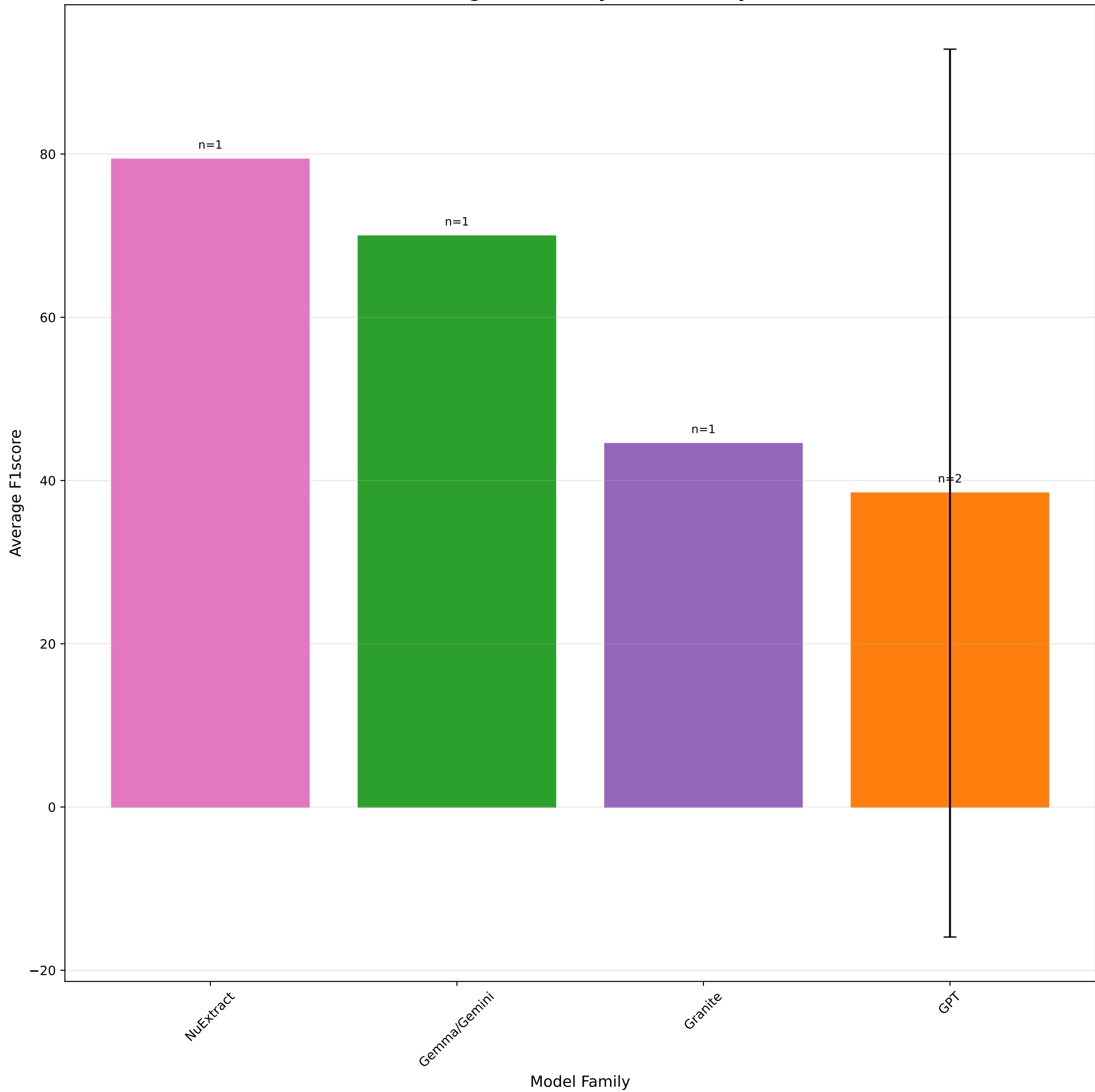
### GROUPED MODEL F1 SCORE STATISTICS:

- Unique Base Models: 5
- Total Test Instances: 54
- Best Performing Model: NuExtract:4B (F1: 79.36)
- Worst Performing Model: granite3.2+LTNER/GPT-NER Prompt (F1: 0.00)
- Overall Average F1: 54.15
- Models with Vision: 3

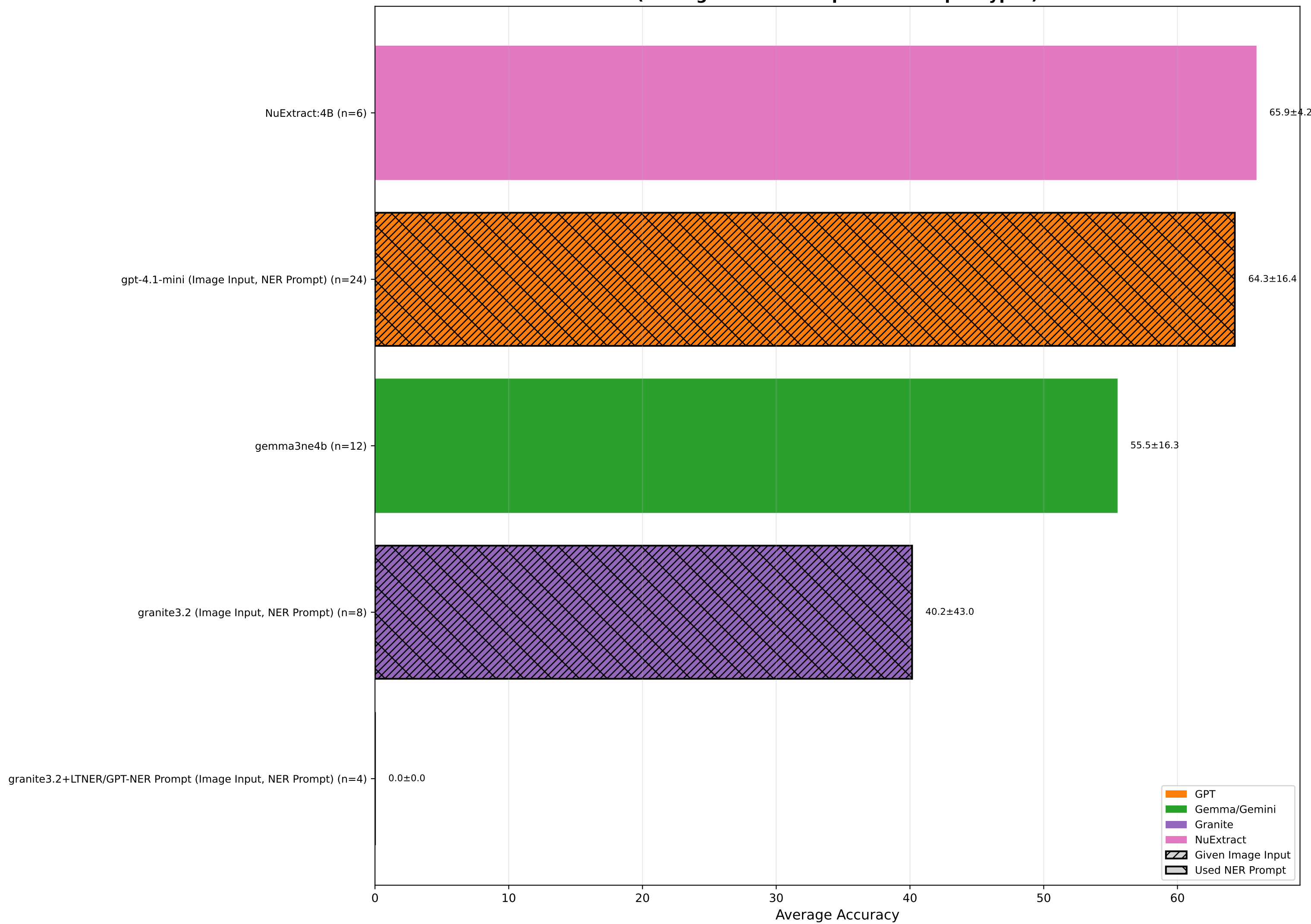
Overall F1score Performance - Models Grouped by Base Name  
(Averaged across hospitals and input types)



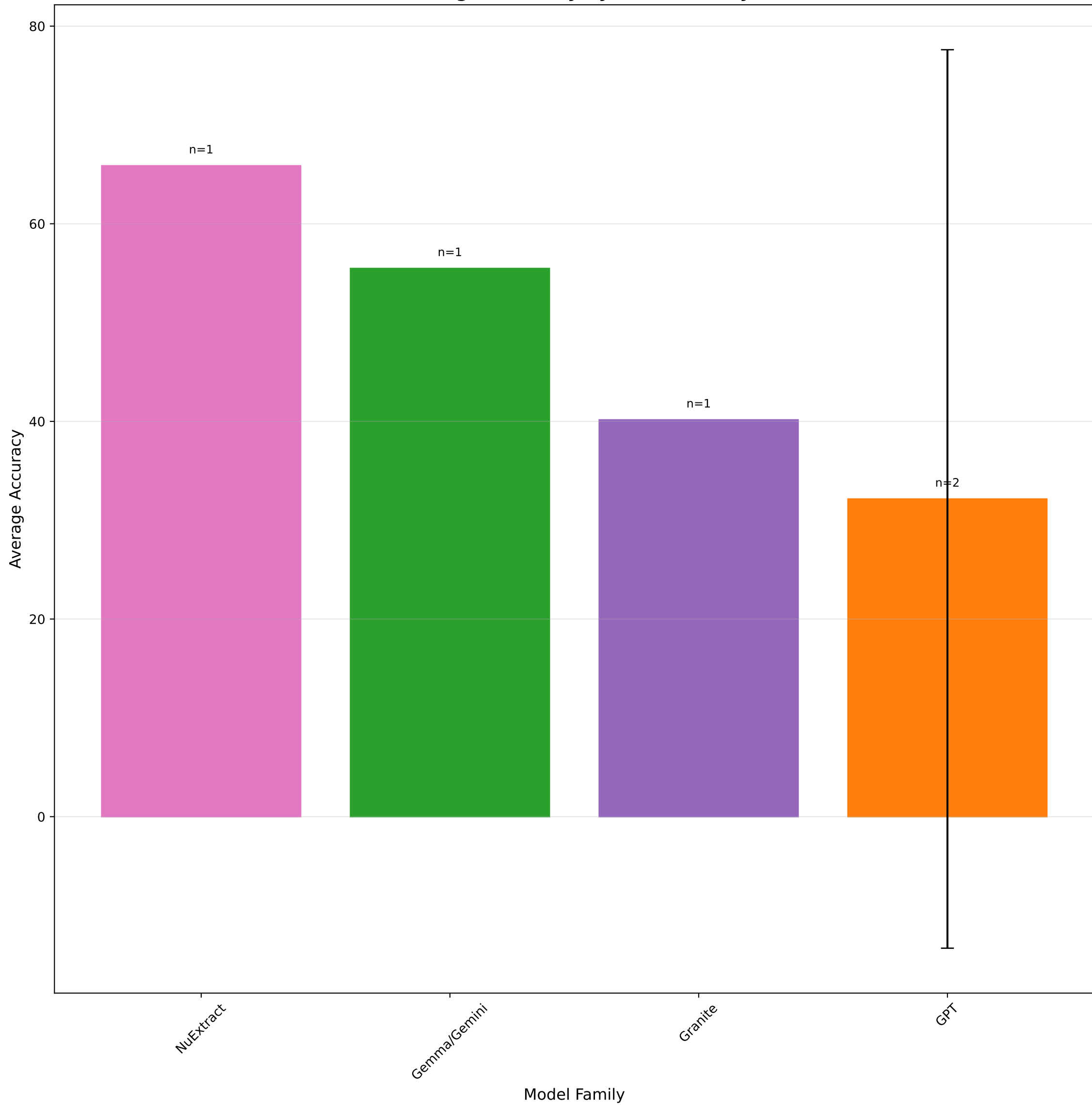
Average F1score by Model Family



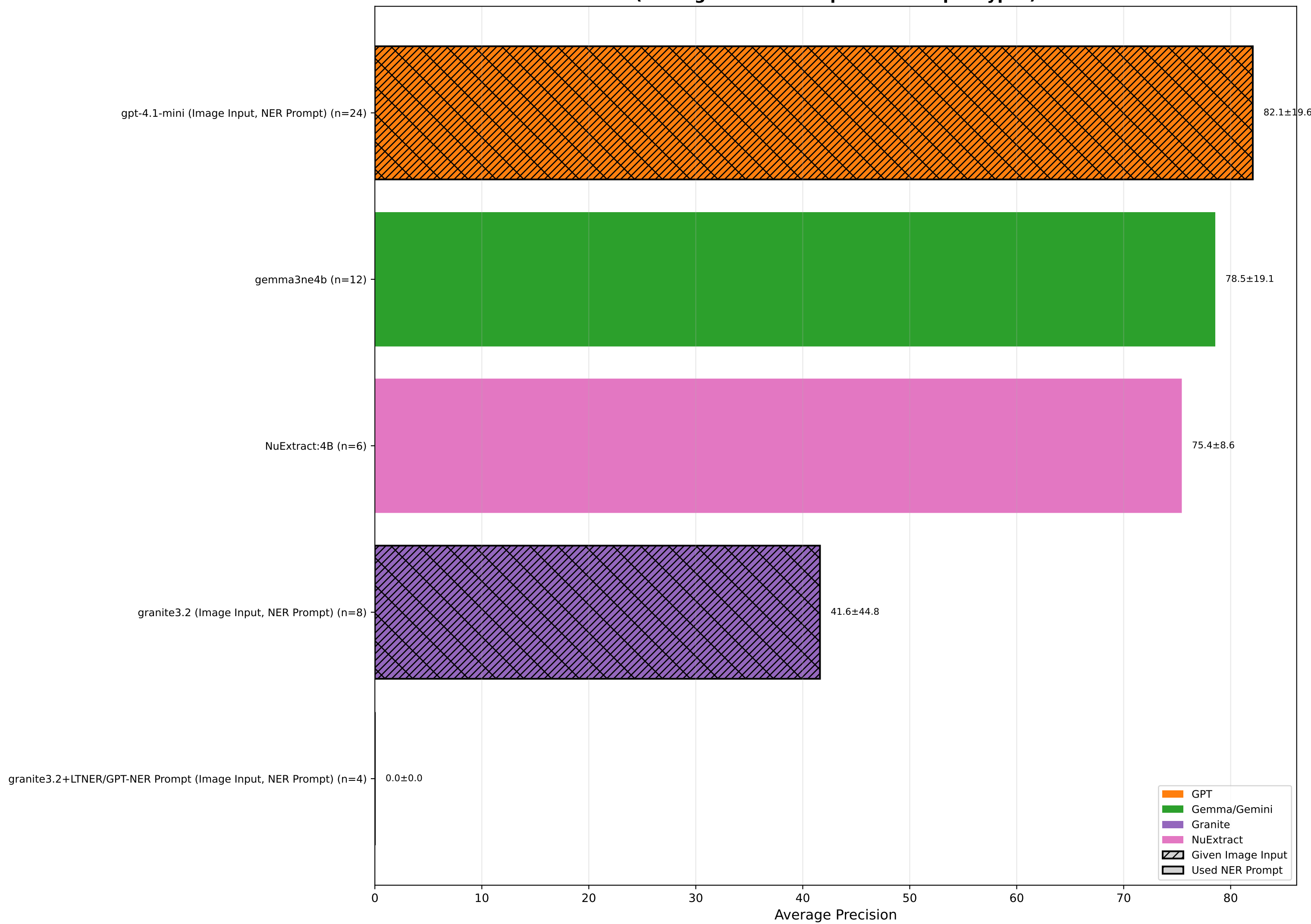
Overall Accuracy Performance - Models Grouped by Base Name  
(Averaged across hospitals and input types)



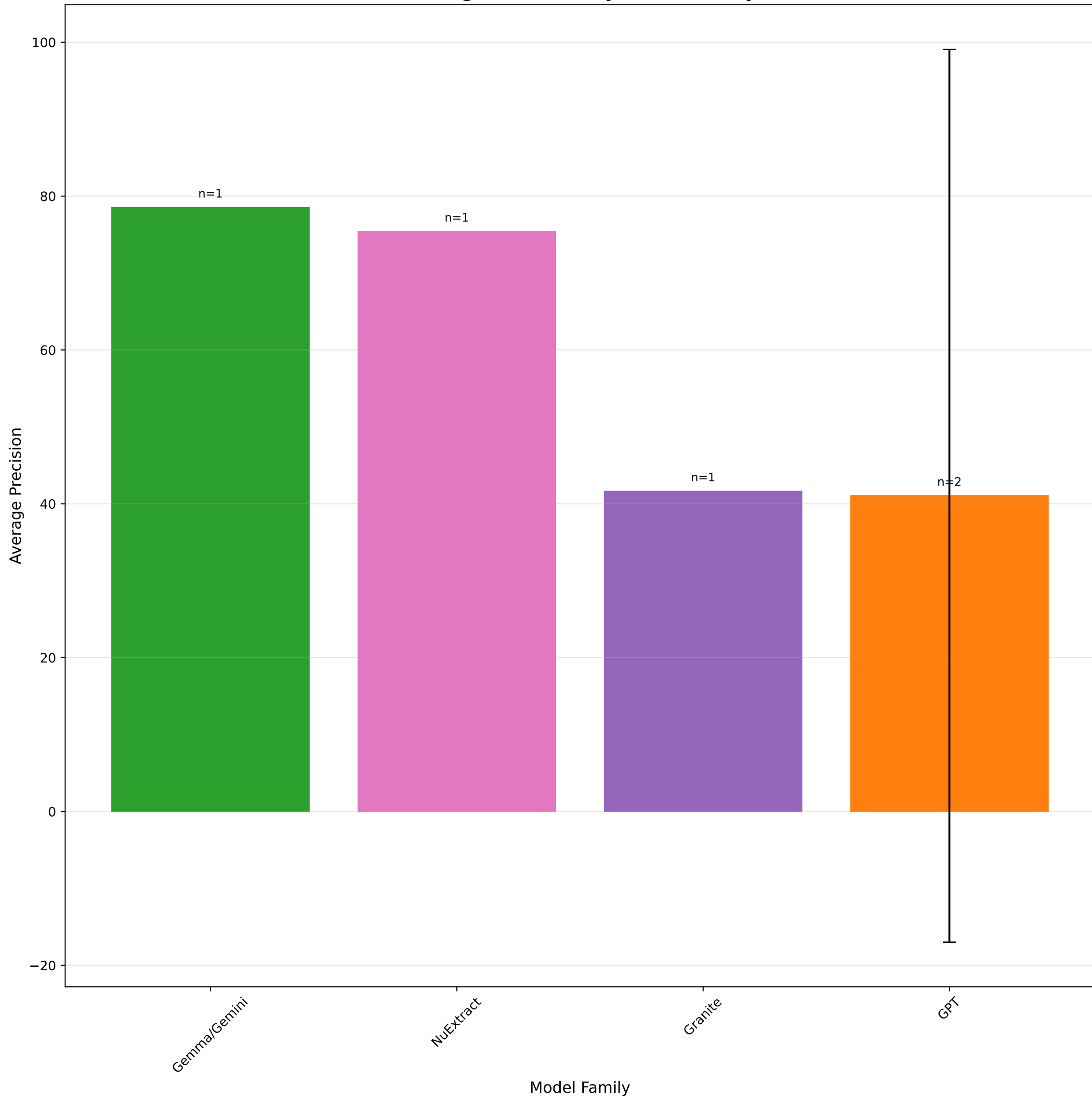
Average Accuracy by Model Family



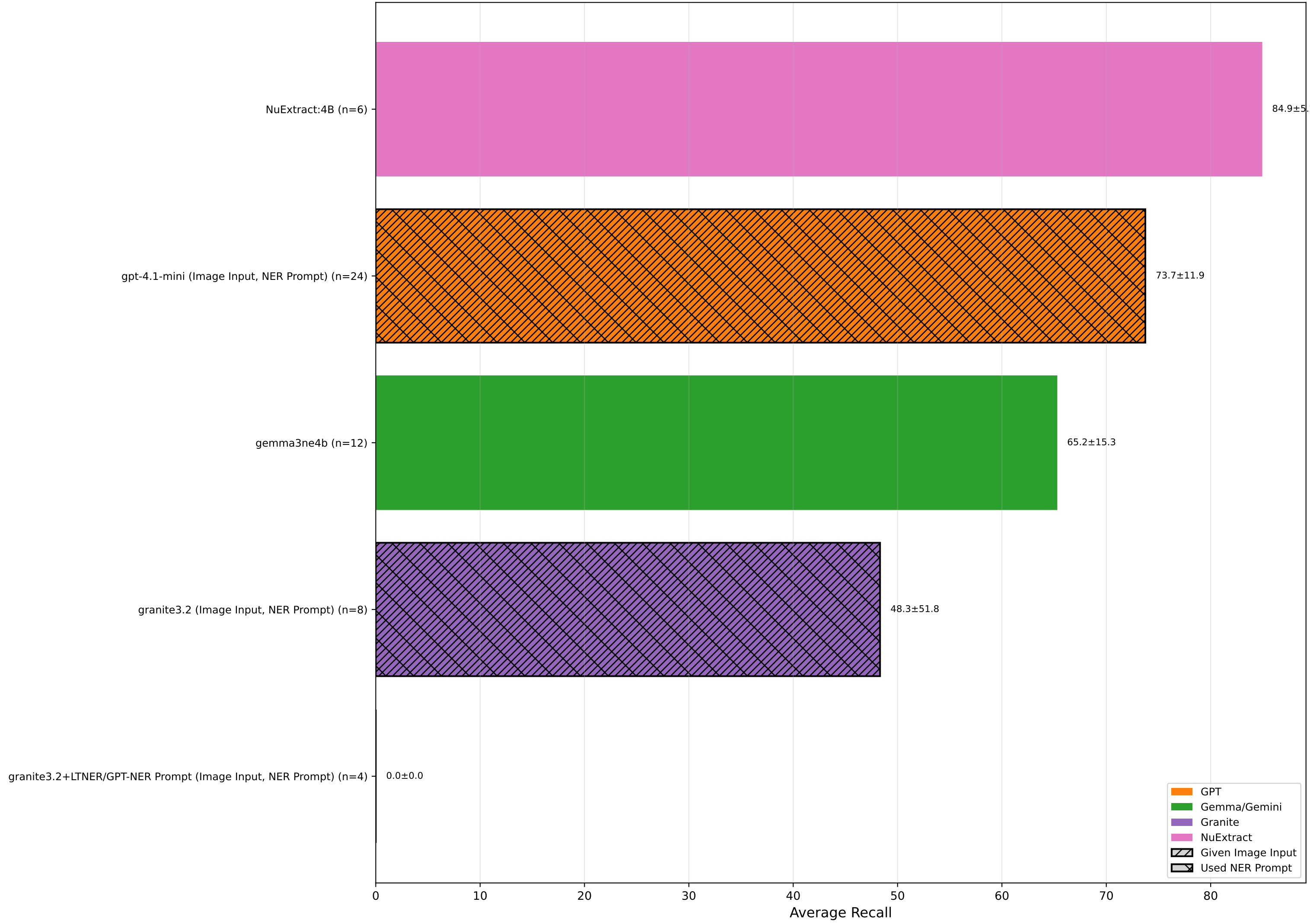
Overall Precision Performance - Models Grouped by Base Name  
(Averaged across hospitals and input types)



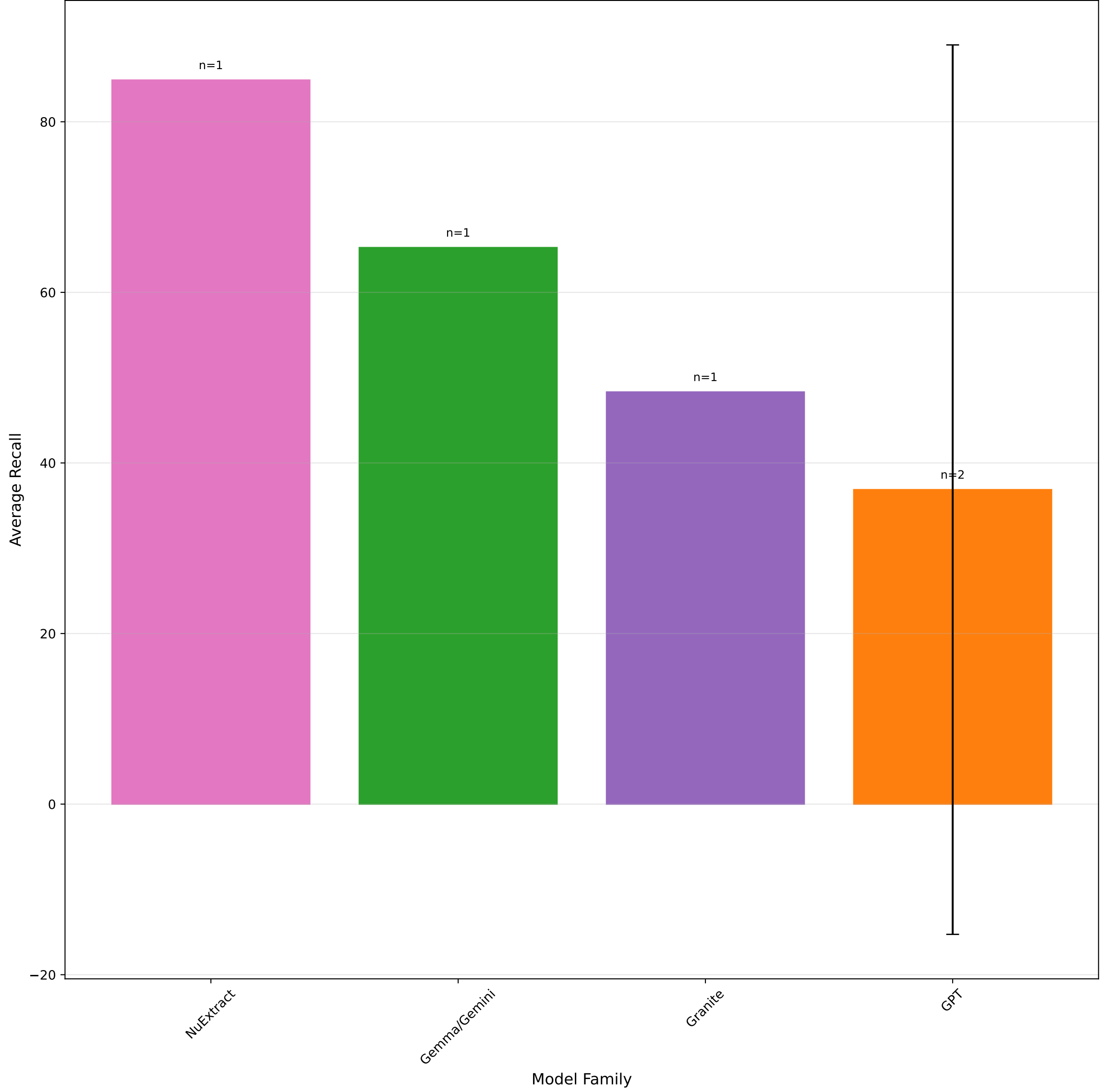
Average Precision by Model Family



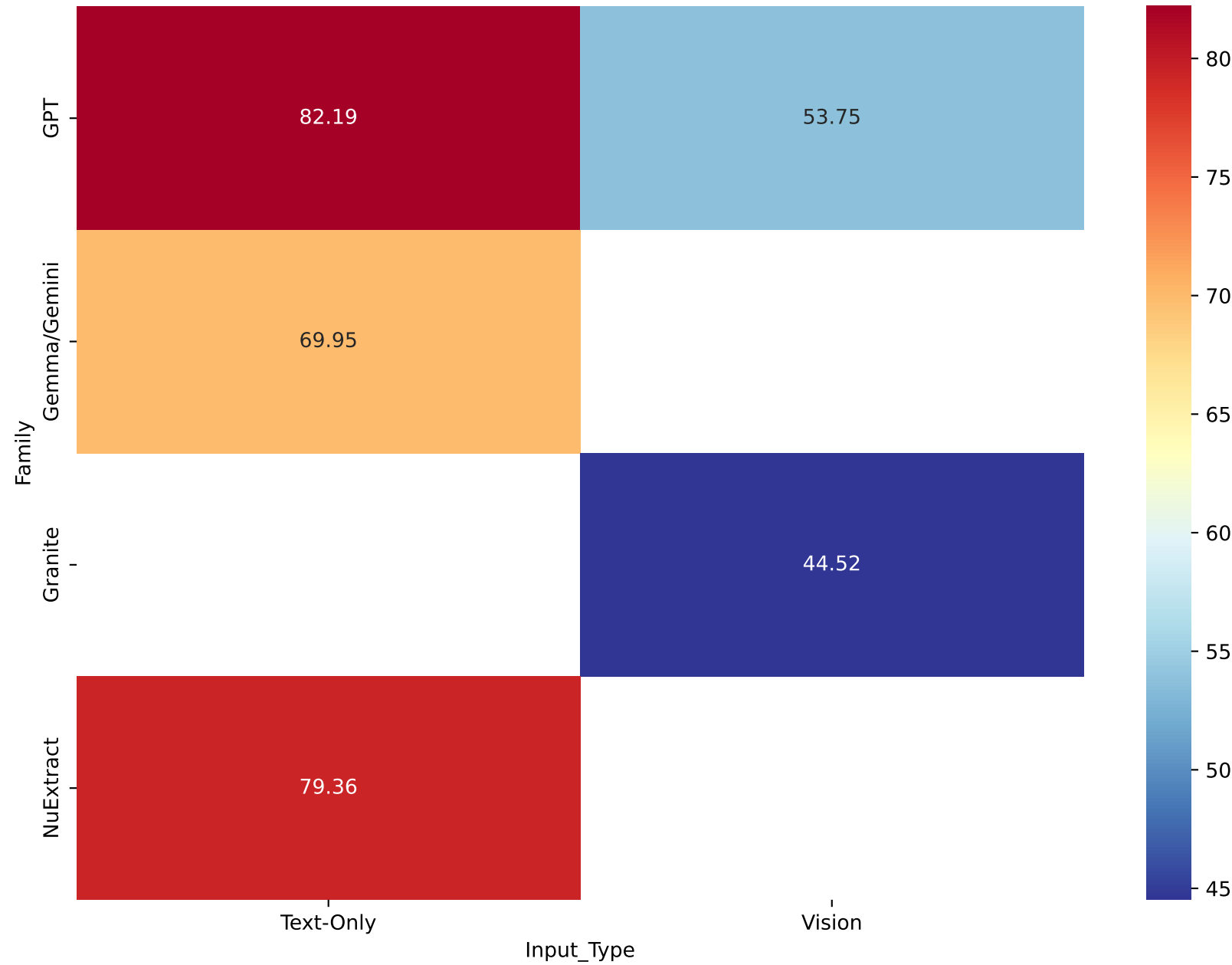
Overall Recall Performance - Models Grouped by Base Name  
(Averaged across hospitals and input types)



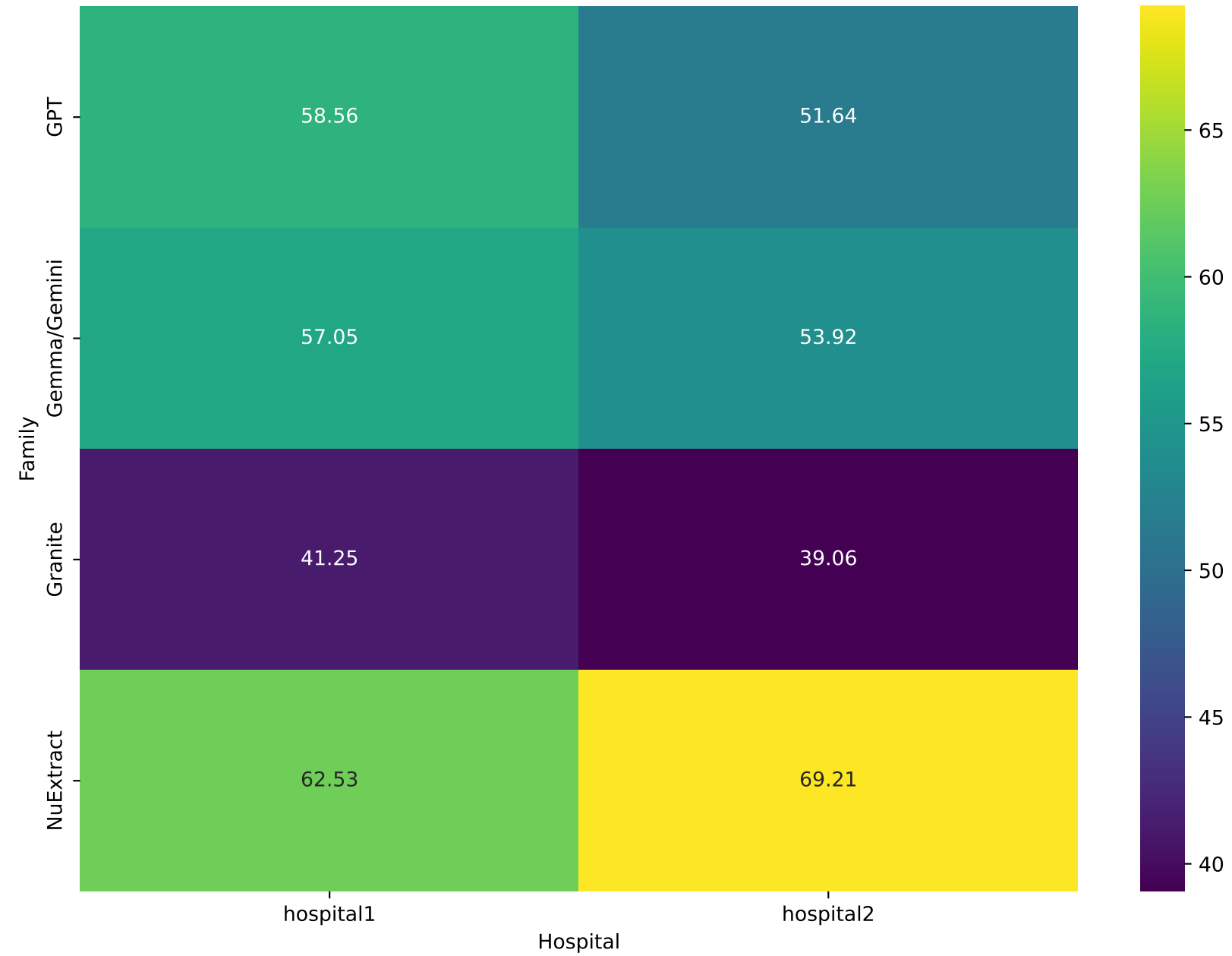
Average Recall by Model Family



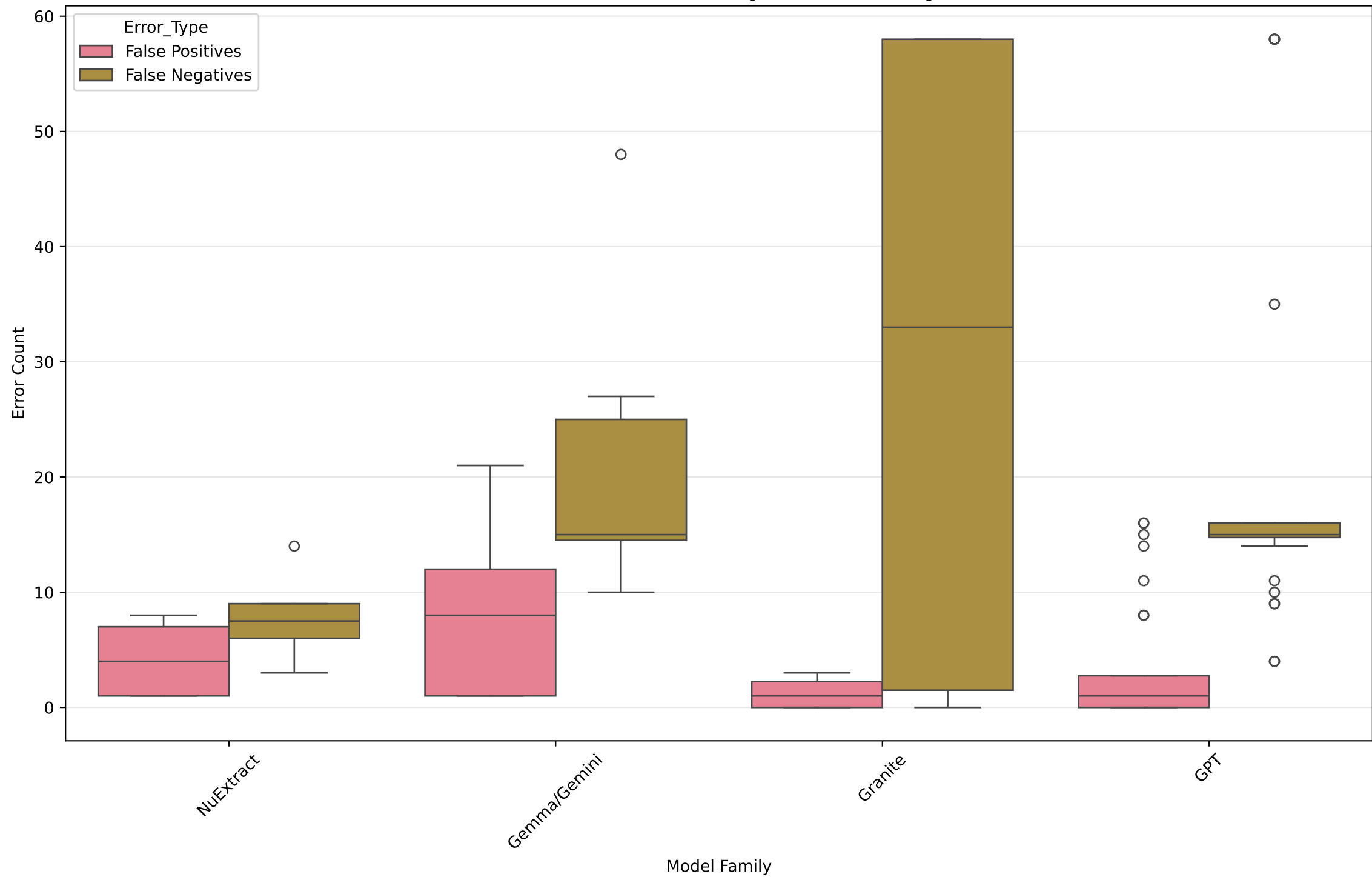
Average F1 Score: Family vs Input Type



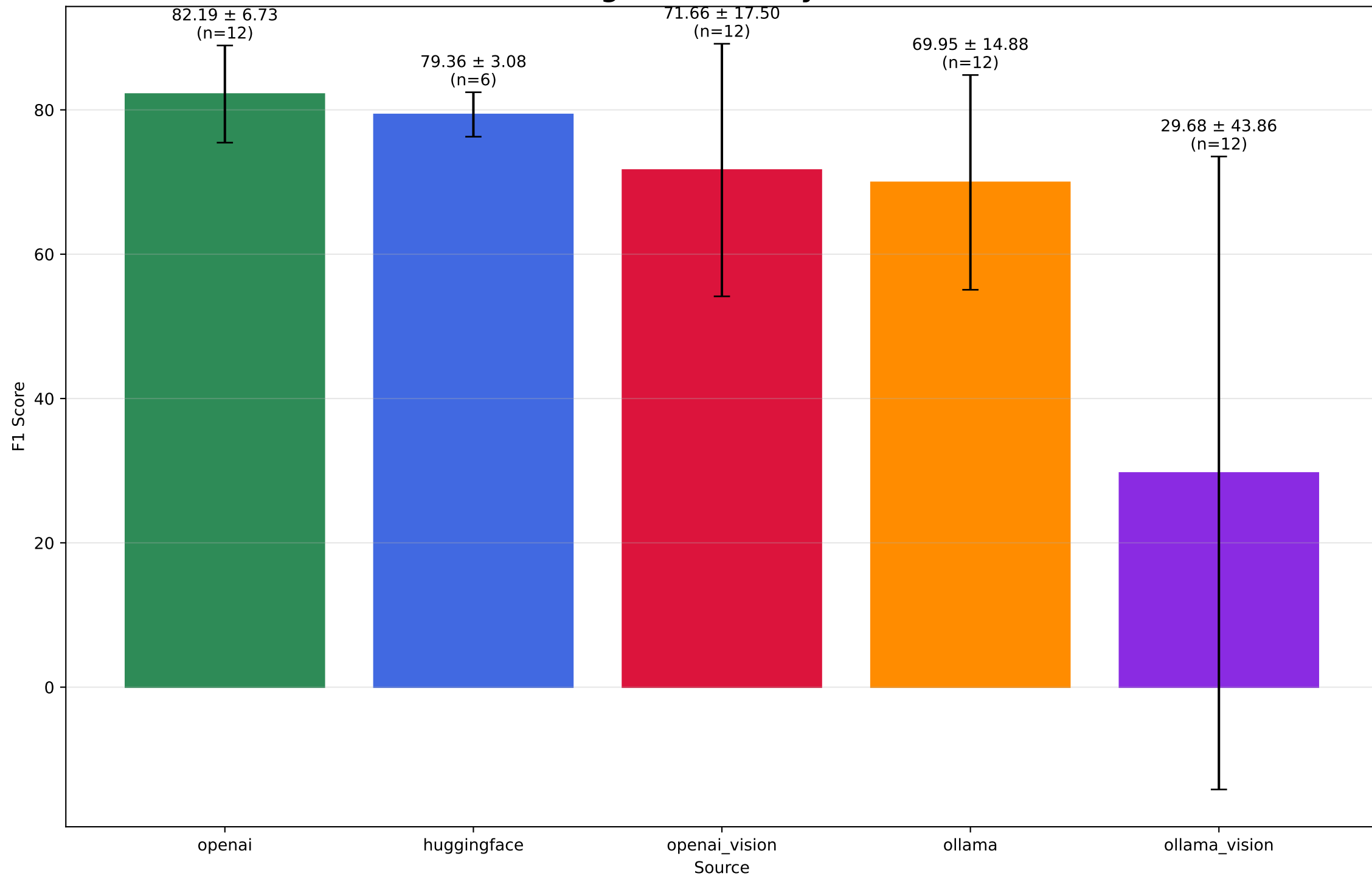
Average Accuracy: Family vs Hospital



**Error Distribution by Model Family**

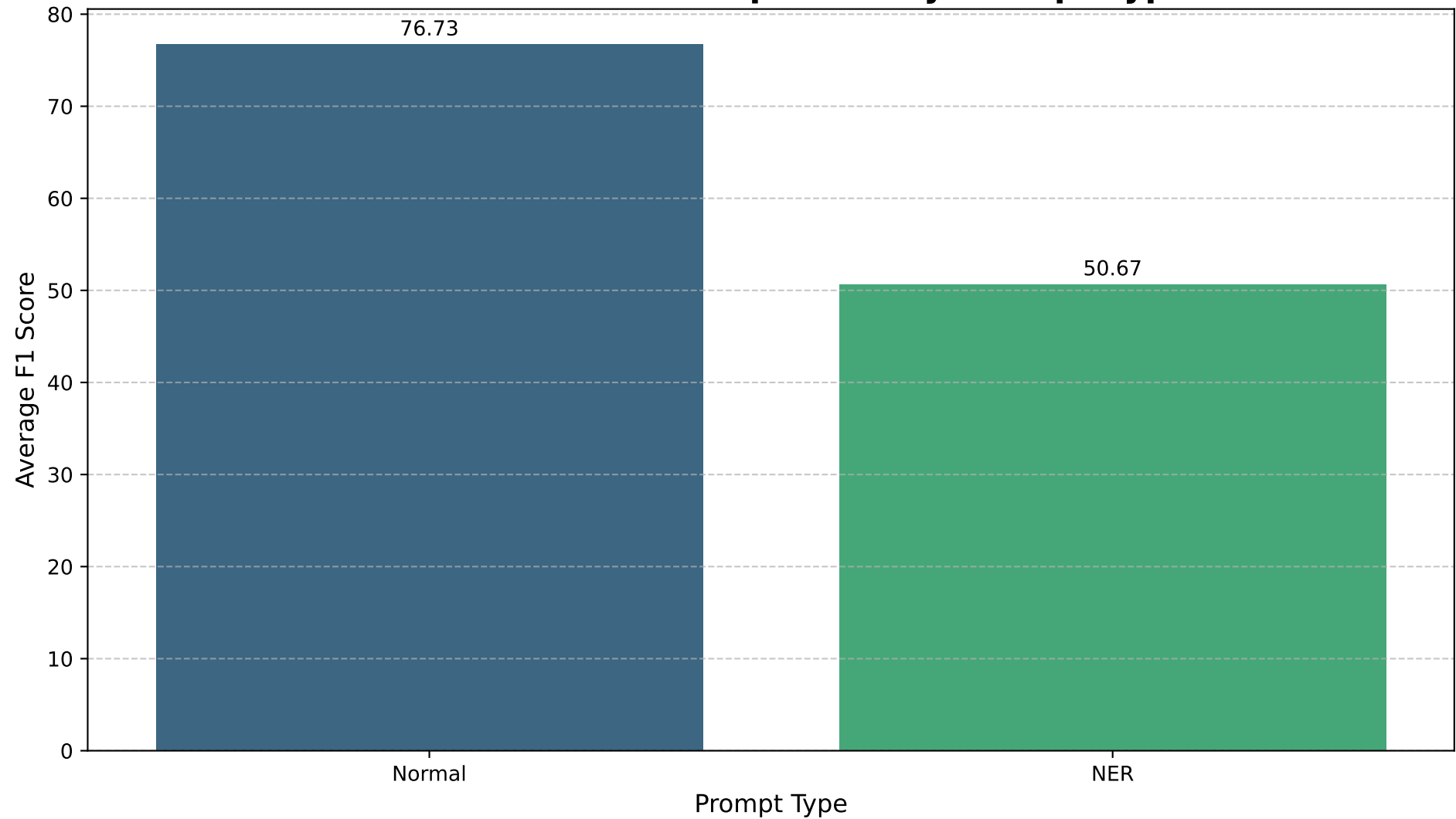


**Average F1 Score by Source**

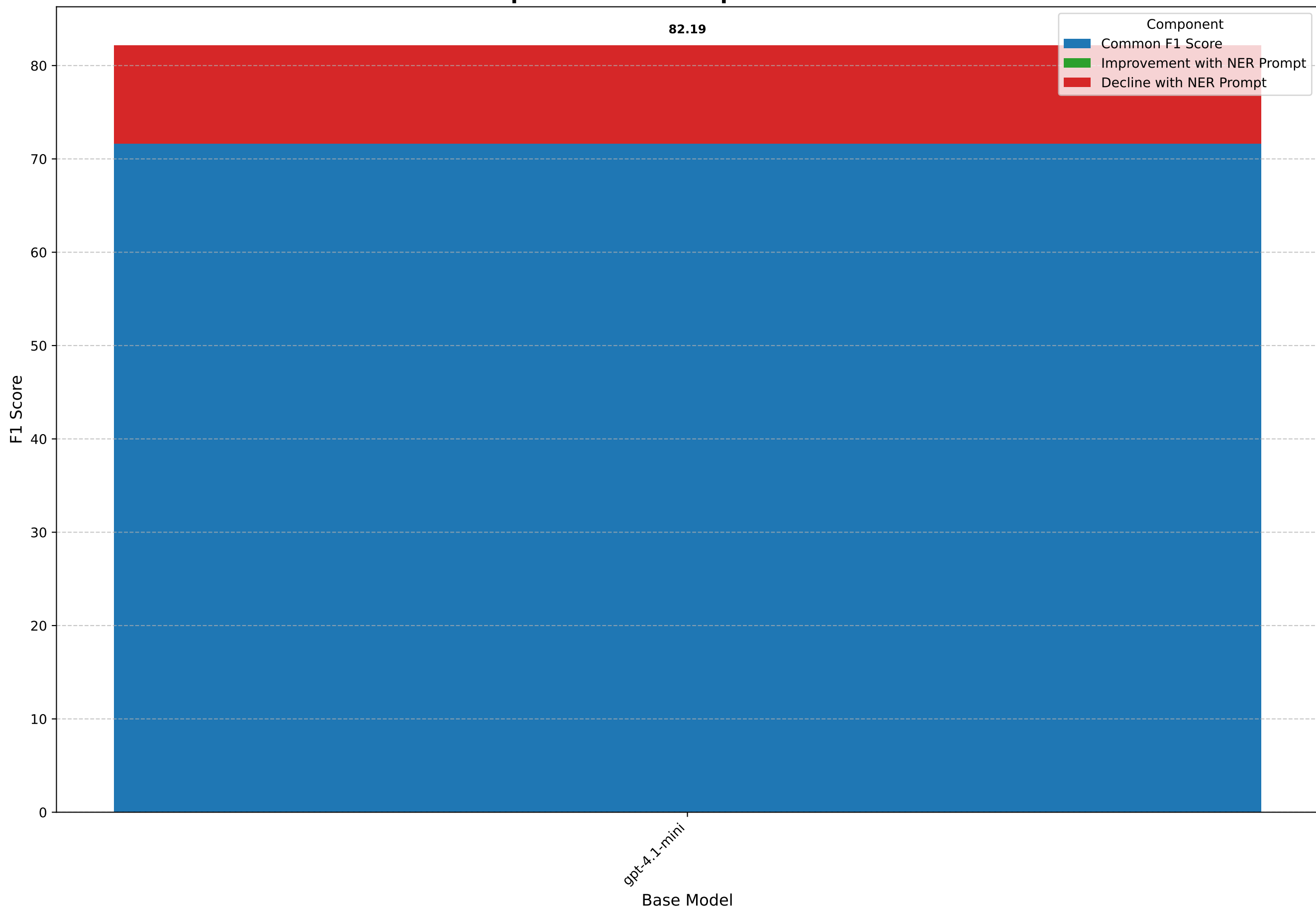




# Overall F1 Score Comparison by Prompt Type



# Impact of NER Prompt on F1 Score



# Impact of Image Input on F1 Score

