

COMPREHENSIVE LLM PERFORMANCE ANALYSIS

=====

Report Generated: 2025-07-07 22:50:26

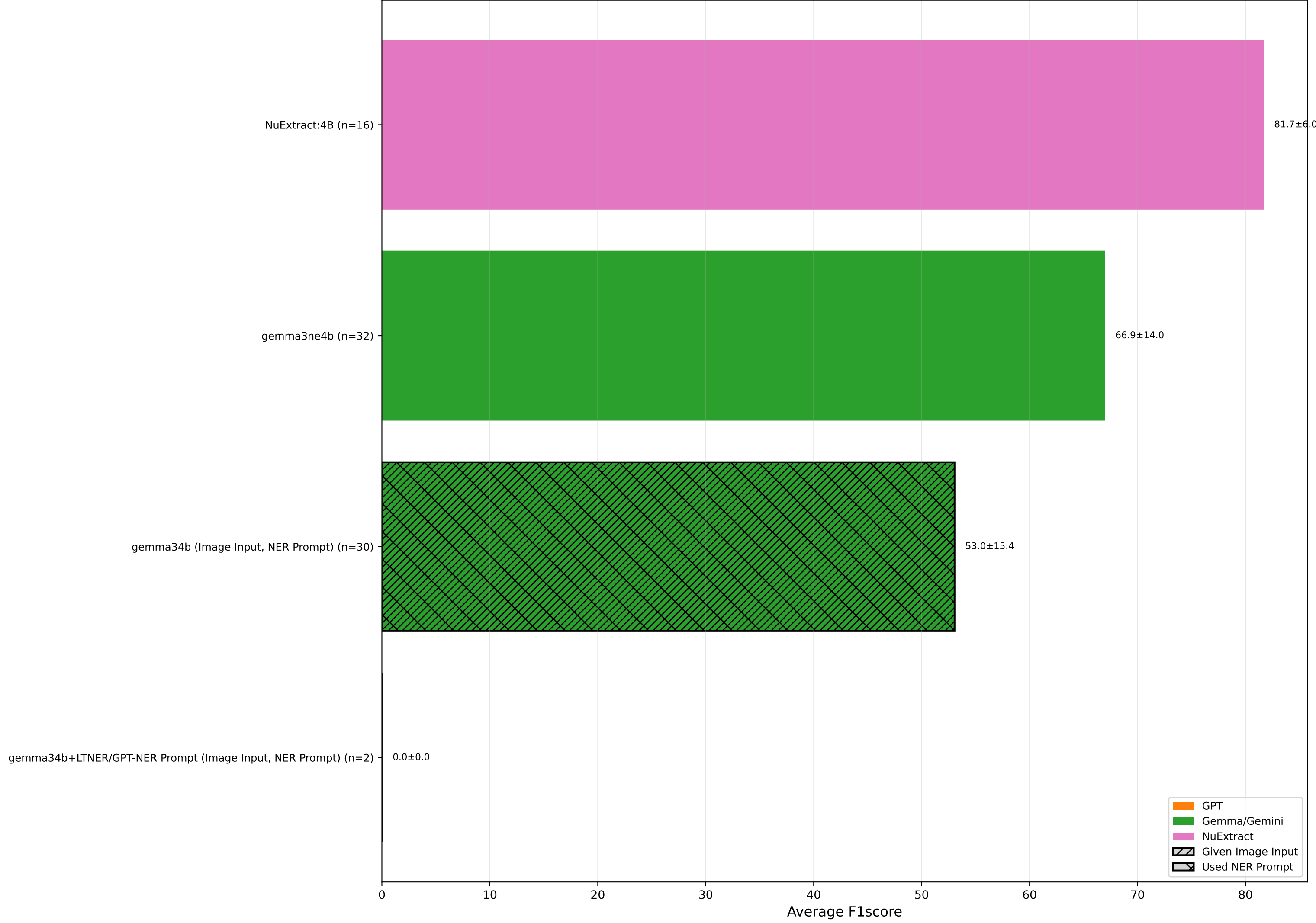
OVERALL PERFORMANCE METRICS:

- Average F1 Score: 62.999 (Std: 19.634)
- Average Accuracy: 48.572 (Std: 18.590)
- Average Precision: 70.781 (Std: 21.963)
- Average Recall: 59.142 (Std: 21.967)

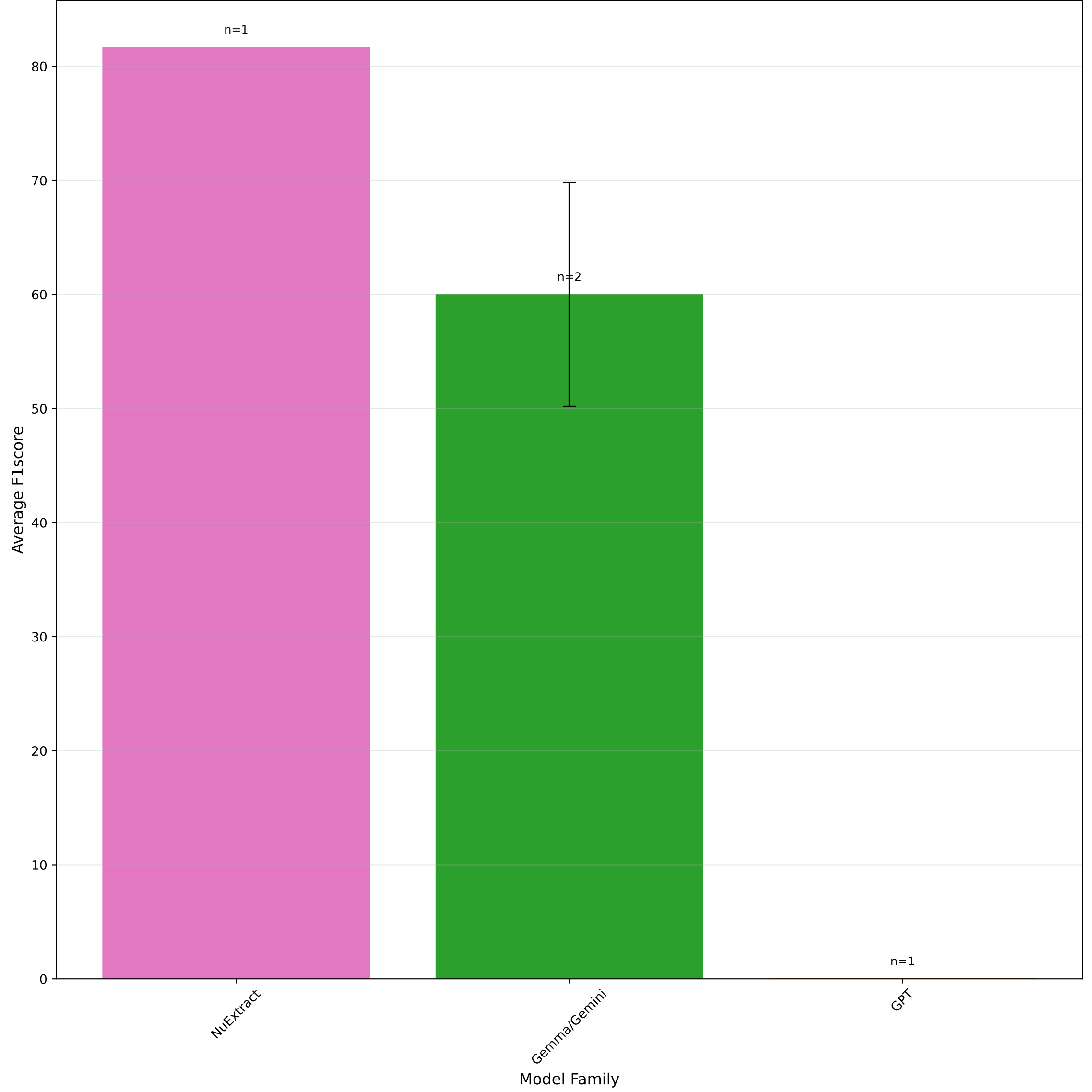
GROUPED MODEL F1 SCORE STATISTICS:

- Unique Base Models: 4
- Total Test Instances: 80
- Best Performing Model: NuExtract:4B (F1: 81.66)
- Worst Performing Model: gemma34b+LTNER/GPT-NER Prompt (F1: 0.00)
- Overall Average F1: 50.41
- Models with Vision: 2

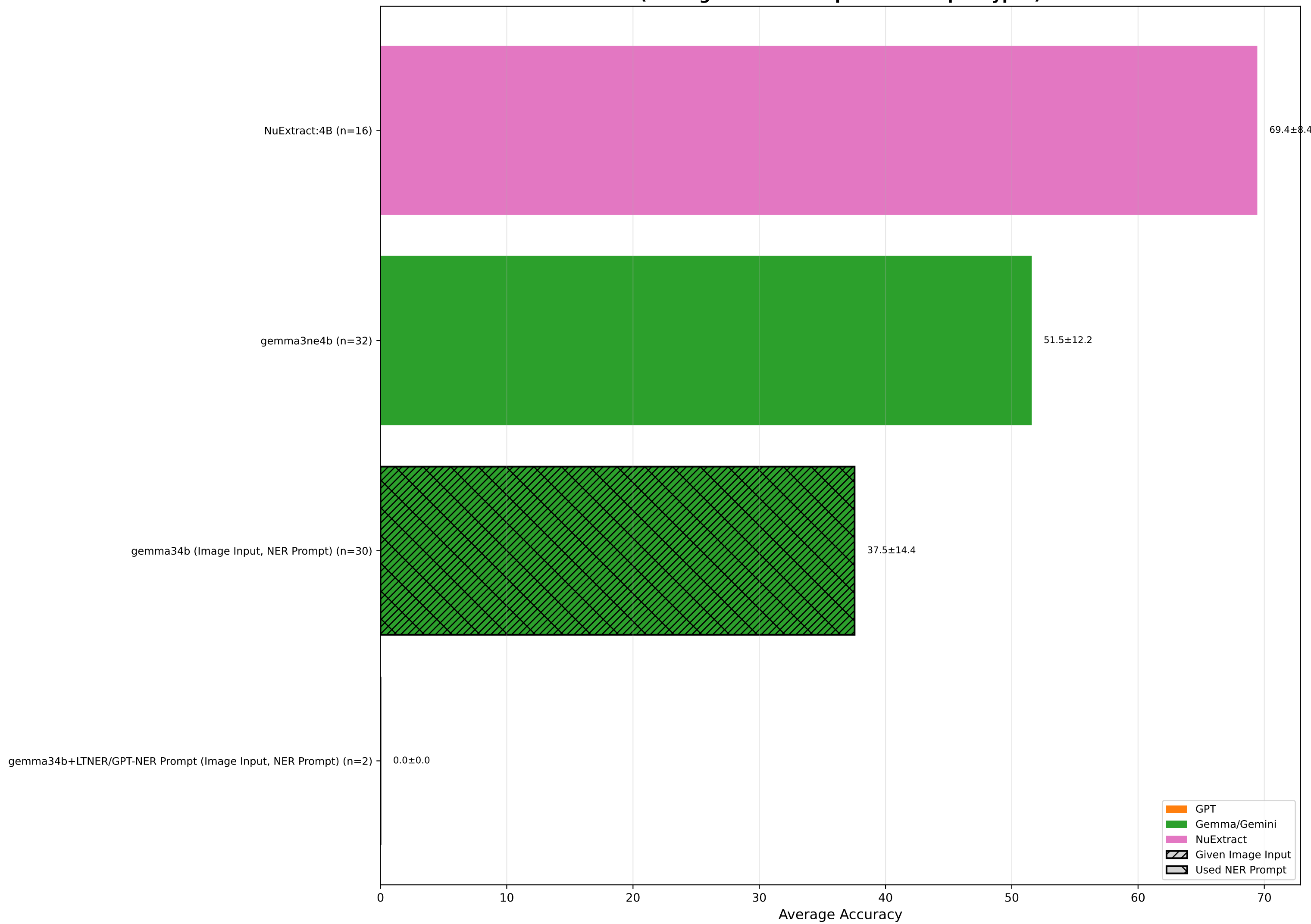
Overall F1score Performance - Models Grouped by Base Name
(Averaged across hospitals and input types)



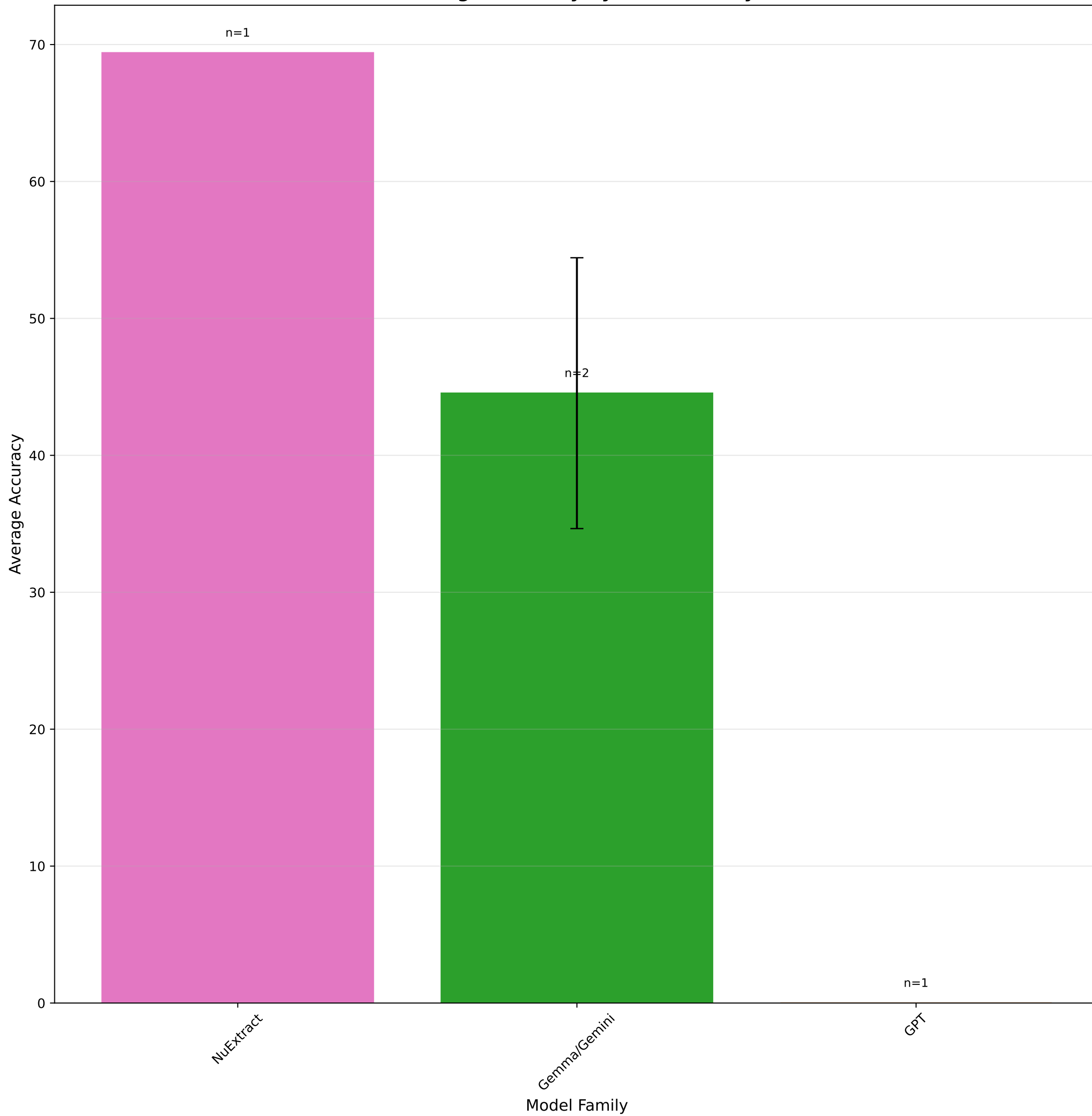
Average F1score by Model Family



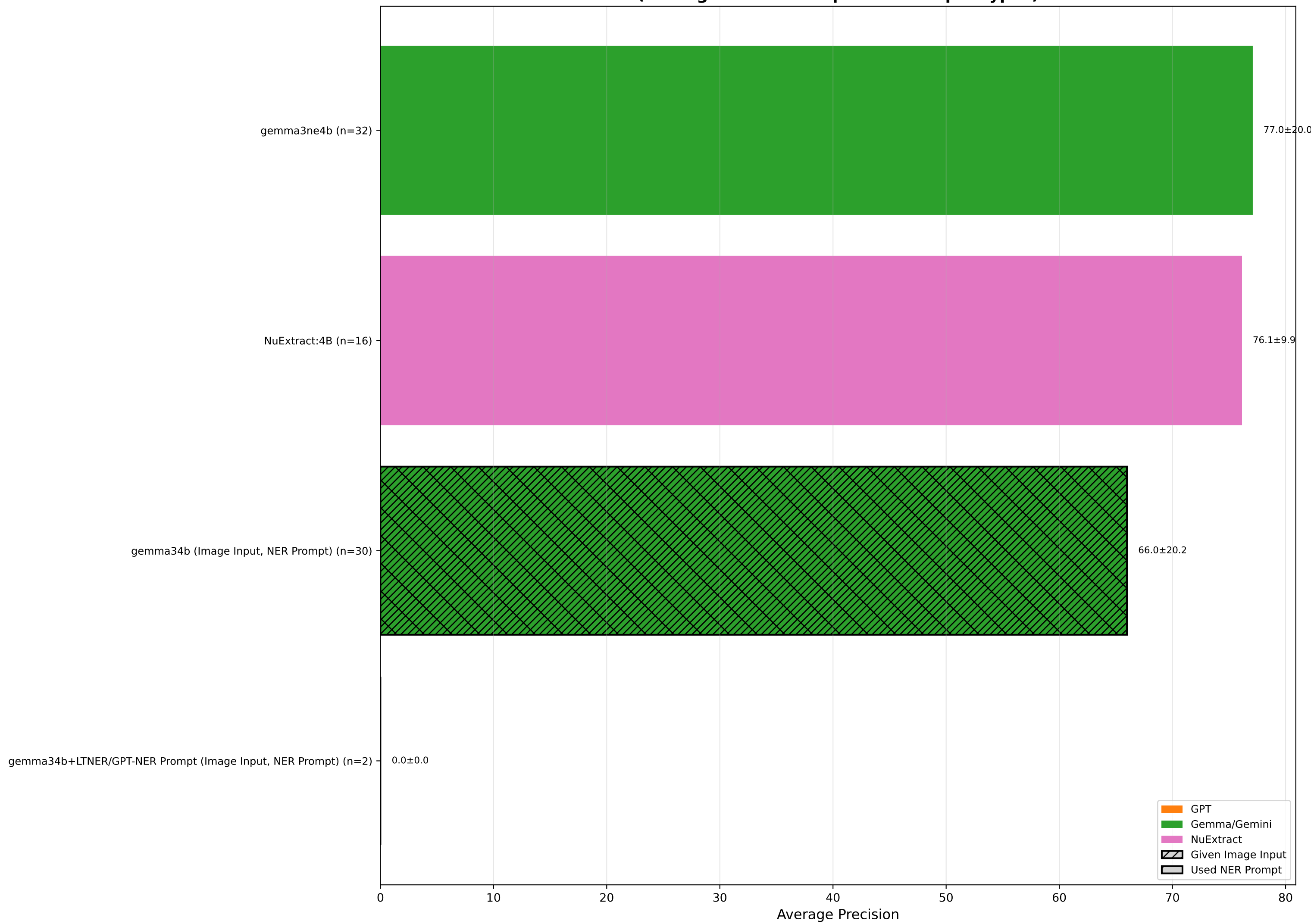
Overall Accuracy Performance - Models Grouped by Base Name
(Averaged across hospitals and input types)



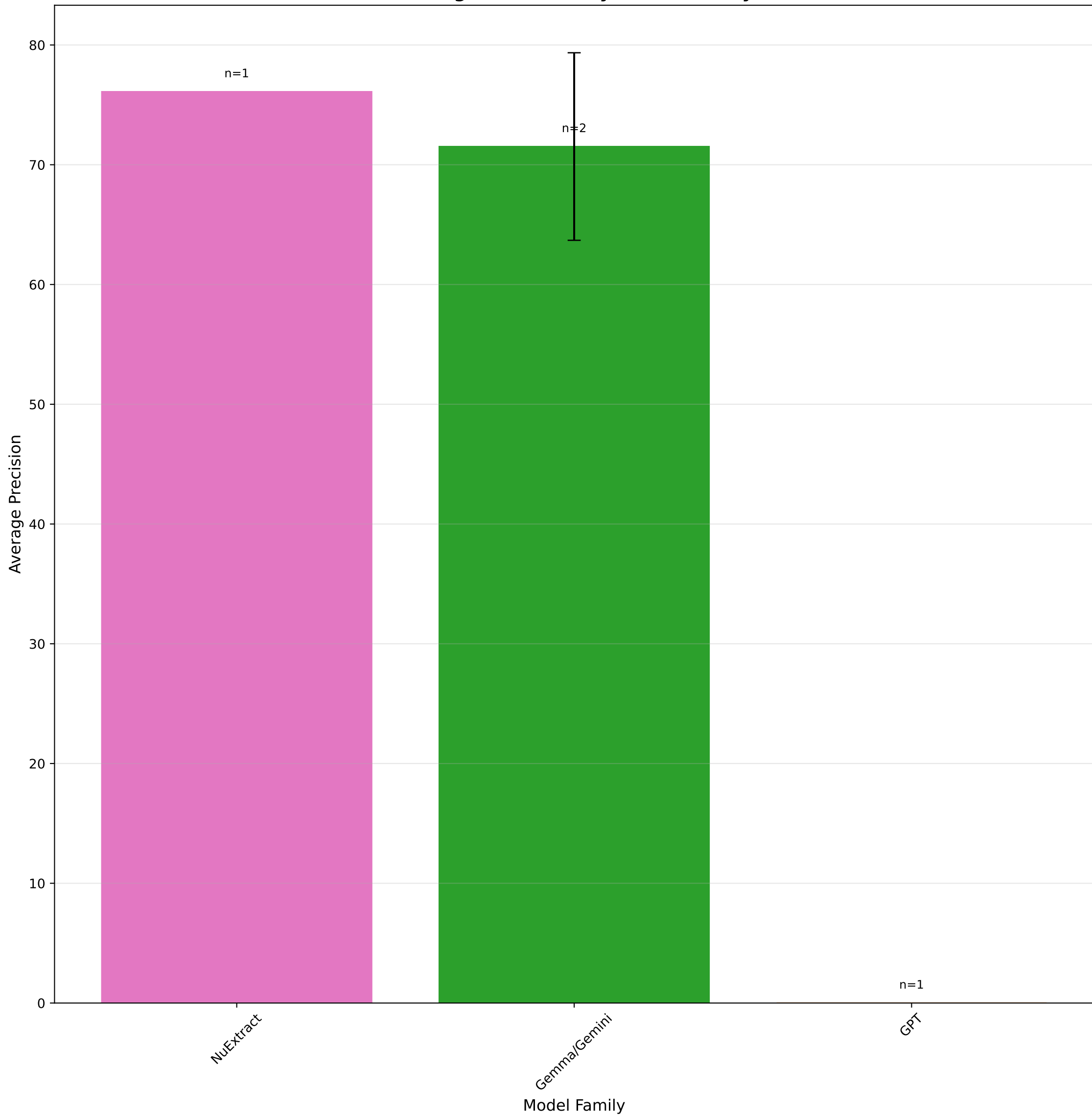
Average Accuracy by Model Family



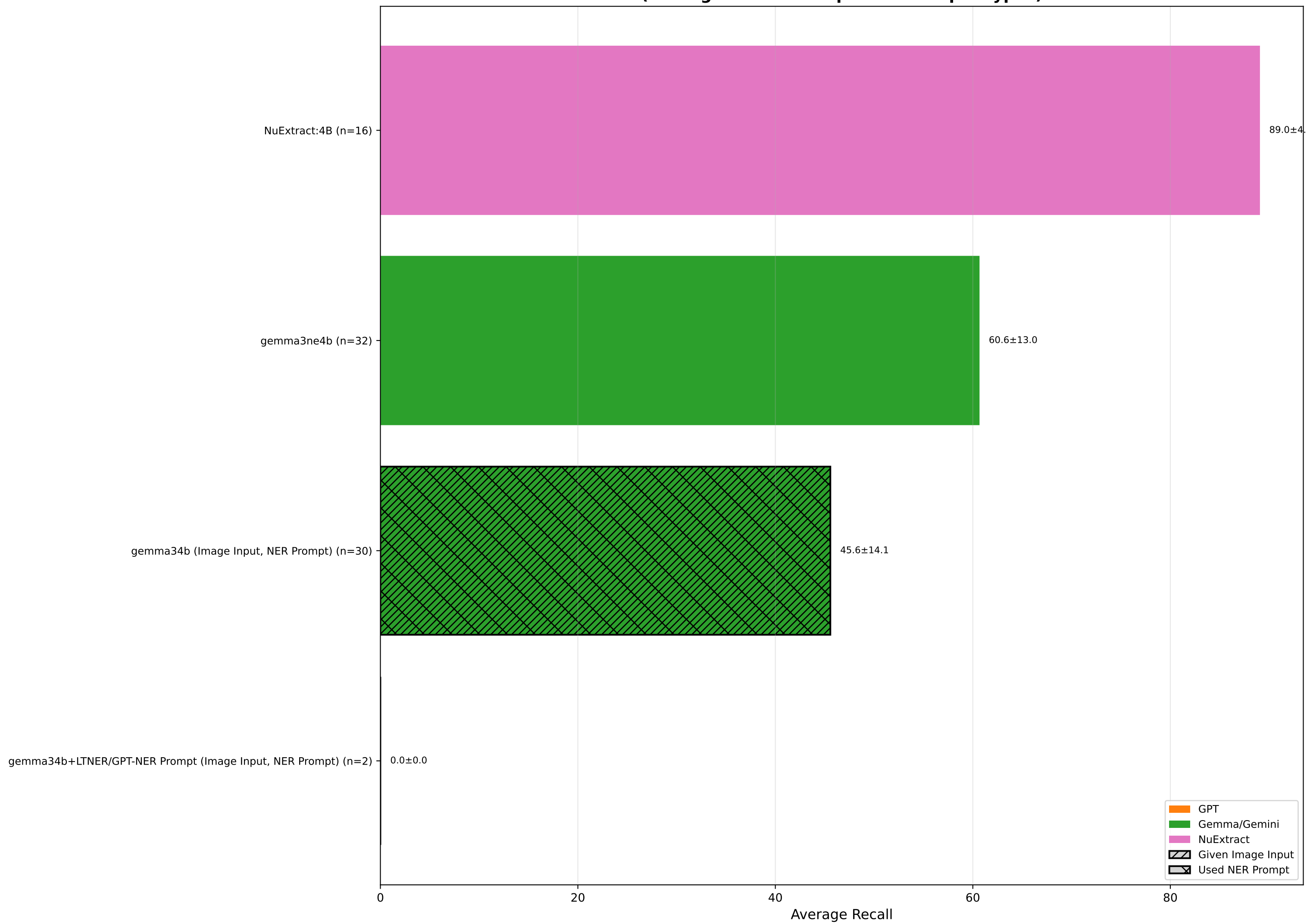
Overall Precision Performance - Models Grouped by Base Name
(Averaged across hospitals and input types)



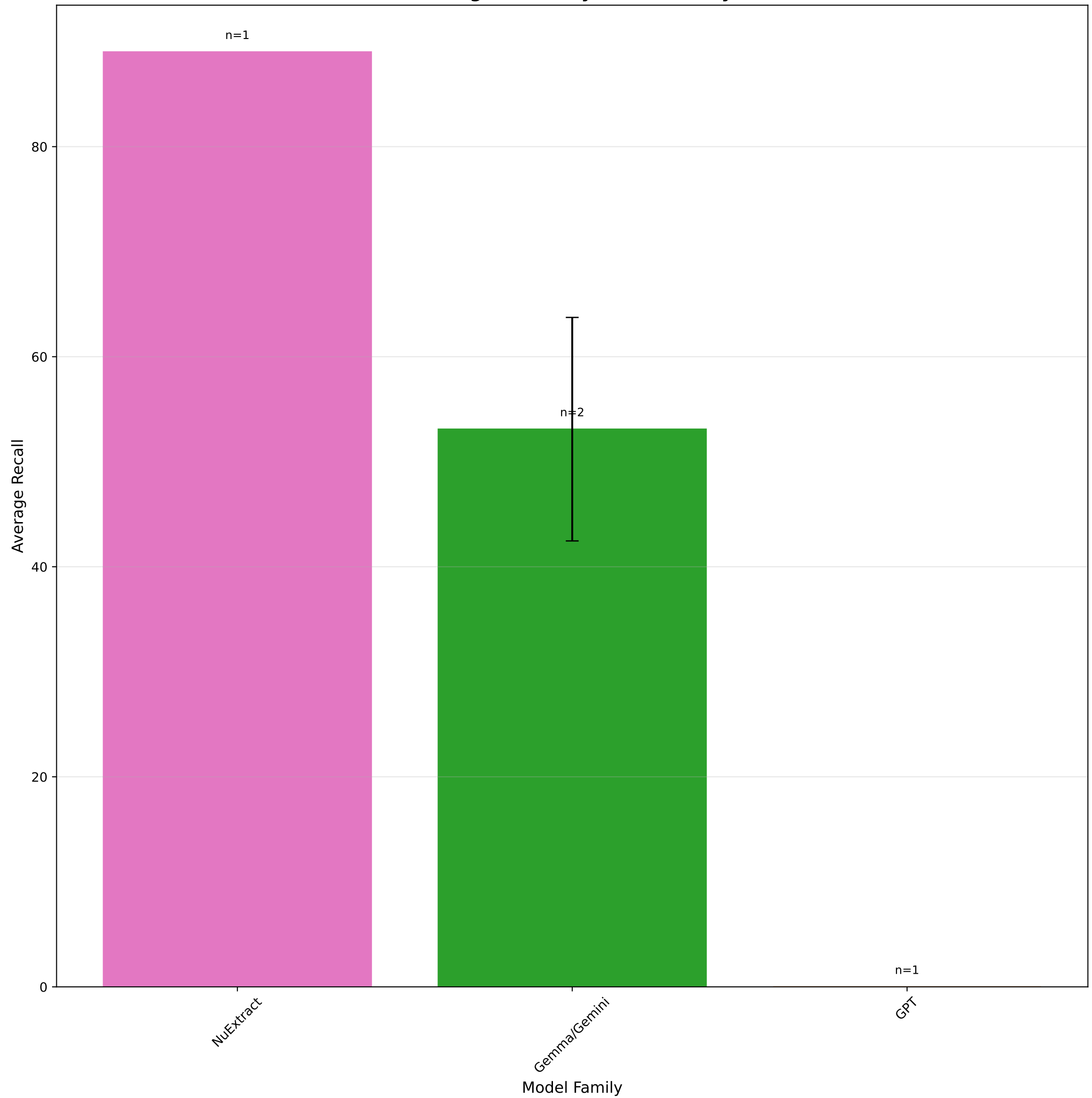
Average Precision by Model Family



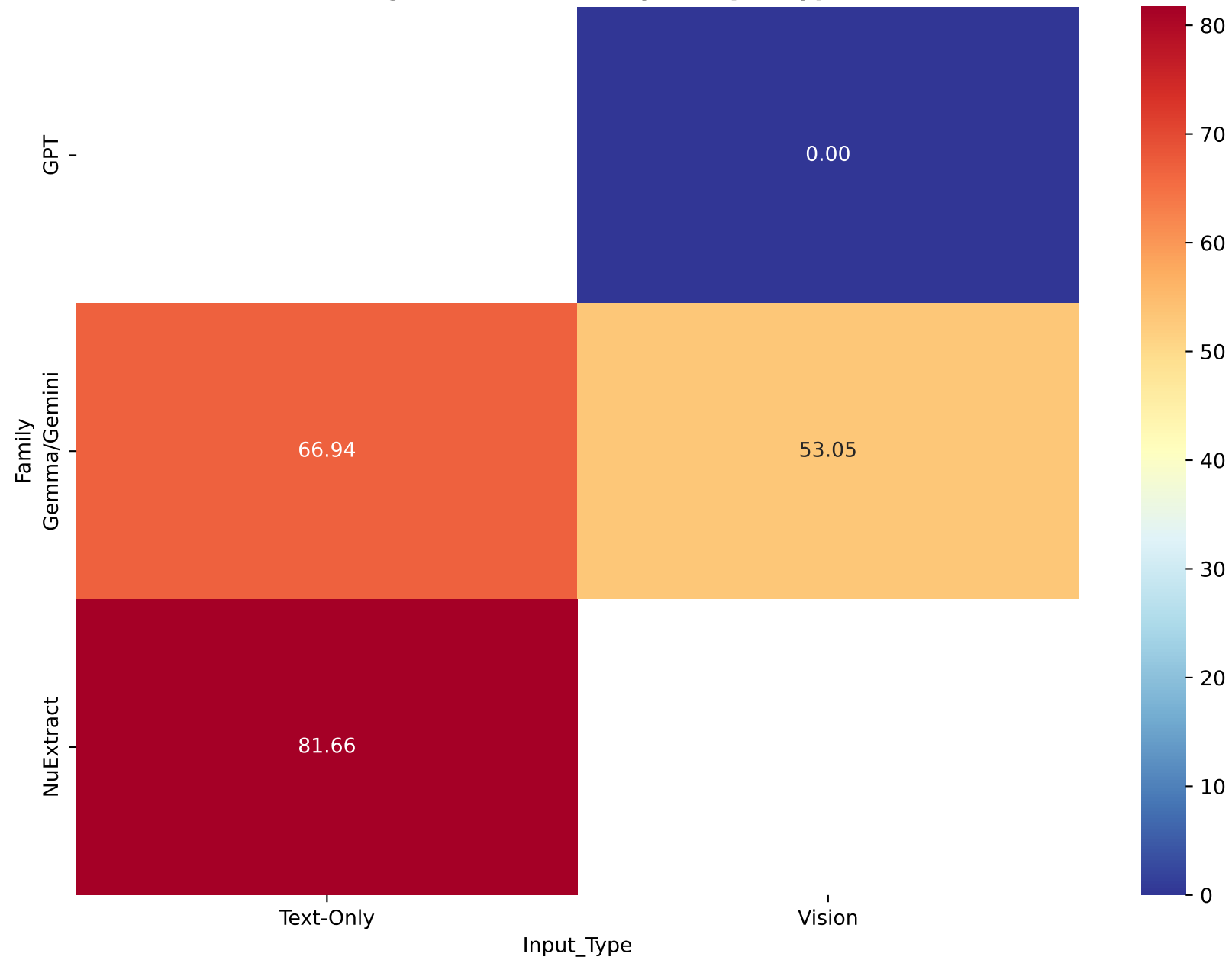
Overall Recall Performance - Models Grouped by Base Name
(Averaged across hospitals and input types)



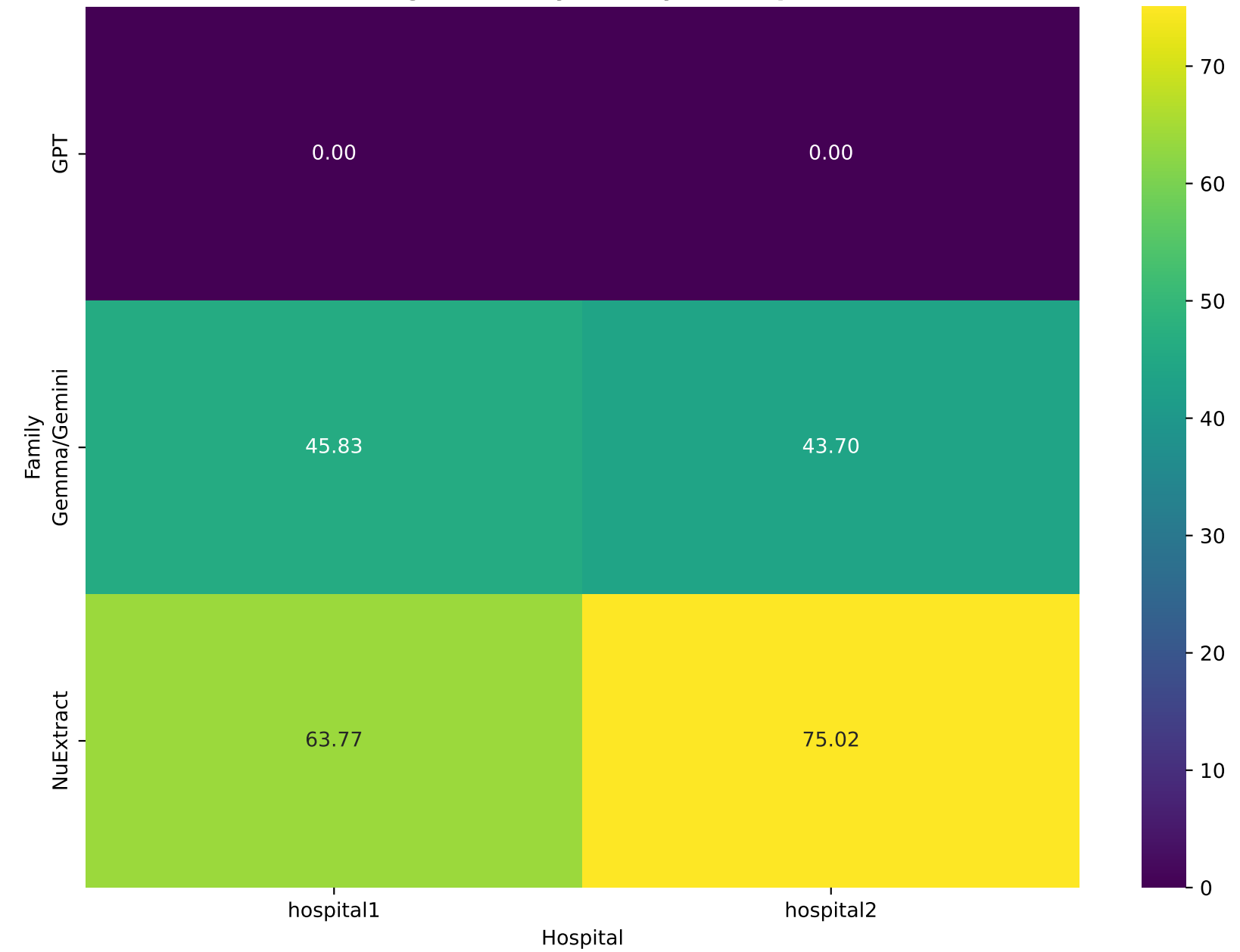
Average Recall by Model Family



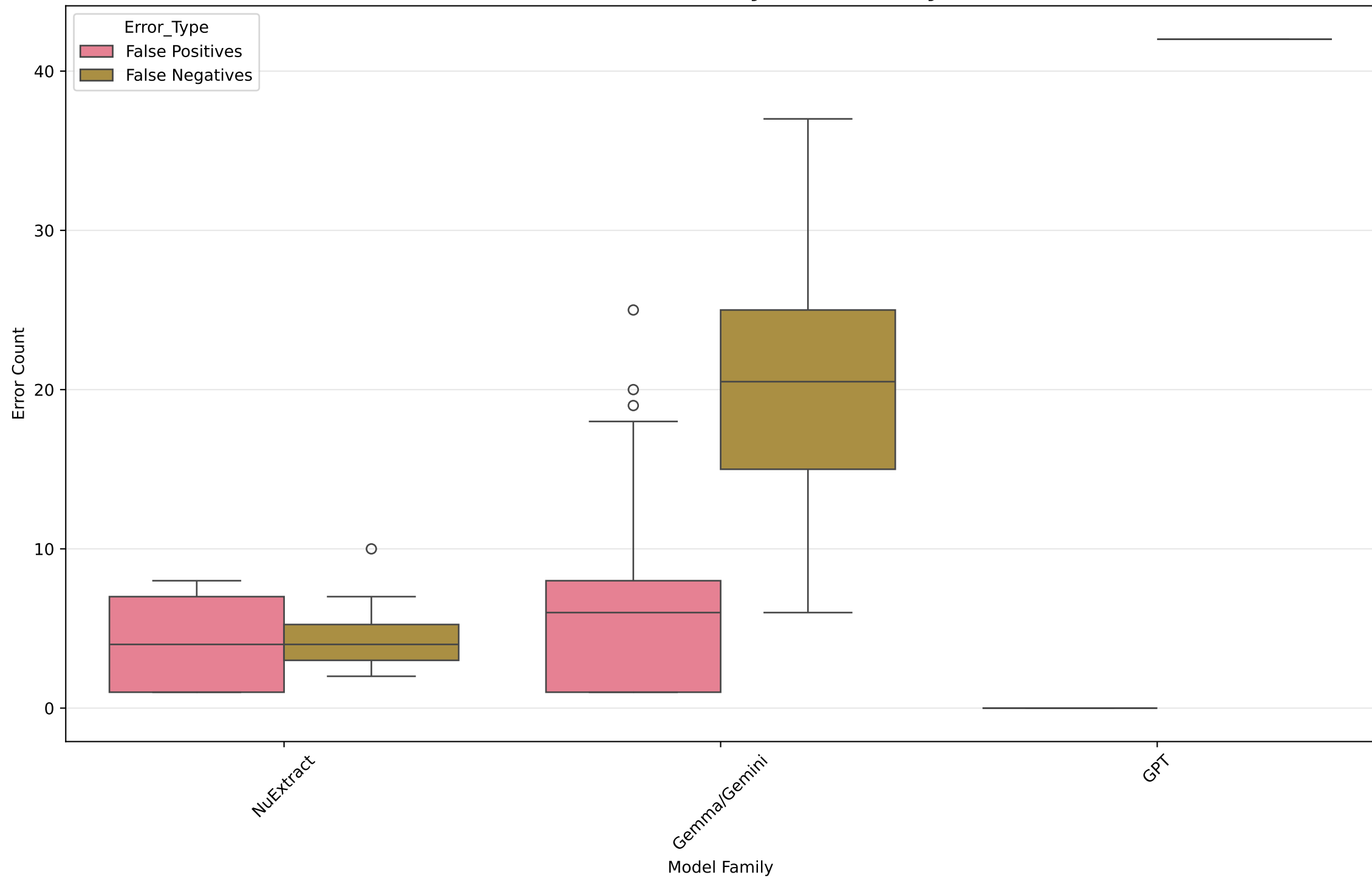
Average F1 Score: Family vs Input Type



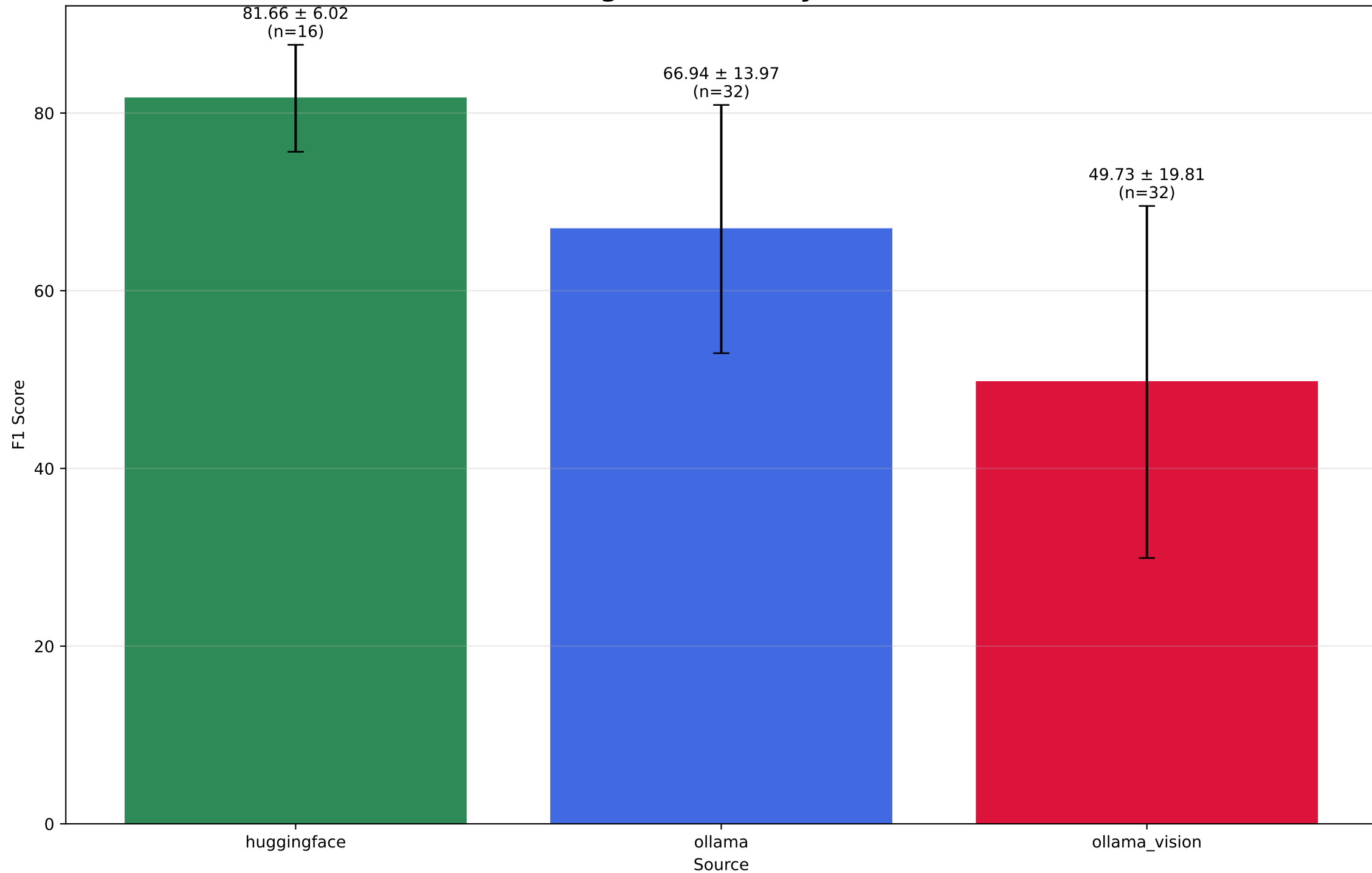
Average Accuracy: Family vs Hospital



Error Distribution by Model Family



Average F1 Score by Source



Overall F1 Score Comparison by Prompt Type

