

COMPREHENSIVE LLM BENCHMARKING ANALYSIS - SUMMARY REPORT

DATASET OVERVIEW:

- Total Records: 145
- Model Families: Qwen, Gemma/Gemini, Granite, Llama, GPT, NuExtract
- Hospitals: hospital1, hospital2
- Vision-enabled Models: 63
- Text-only Models: 82
- Ollama Models: 91
- Commercial Models: 54

TOP PERFORMERS:

- qwen2.5vl3b\*ImageInput\* (Qwen, 3.0B, Vision)  
F1: 98.86, Accuracy: 98.86
- NuExtract:0.5B (NuExtract, 0.5B, Text-only)  
F1: 97.62, Accuracy: 97.62
- NuExtract:2B (NuExtract, 2.0B, Text-only)  
F1: 97.62, Accuracy: 97.62
- NuExtract:2B (NuExtract, 2.0B, Text-only)  
F1: 97.62, Accuracy: 97.62
- NuExtract:0.5B (NuExtract, 0.5B, Text-only)  
F1: 97.62, Accuracy: 97.62

FAMILY PERFORMANCE RANKING (by F1 Score):

- NuExtract: 97.62 (n=8.0)
- GPT: 59.17 (n=46.0)
- Gemma/Gemini: 58.13 (n=29.0)
- Llama: 54.88 (n=15.0)
- Qwen: 50.90 (n=28.0)
- Granite: 47.82 (n=19.0)

KEY INSIGHTS:

- Vision Models Avg F1: 53.33
- Text-only Models Avg F1: 60.80
- Hospital 1 Avg F1: 57.70
- Hospital 2 Avg F1: 57.40
- Ollama Models Avg F1: 53.22
- Commercial Models Avg F1: 64.87

Grouped Model F1 Score Statistics:

- Unique Base Models: 13
- Total Test Instances: 145
- Best Performing Model: NuExtract:2B (F1: 97.62)
- Worst Performing Model: granite3.32b (F1: 37.28)
- Overall Average F1: 60.04
- Models with Vision: 6

Top 5 Performers:

- gpt-4.1-nano (GPT, with Vision): F1 = 62.09 ± 11.01
- llama3.23b (Llama, Text-only): F1 = 62.60 ± 11.41
- gemma31b (Gemma/Gemini, Text-only): F1 = 64.94 ± 11.08
- NuExtract:0.5B (NuExtract, Text-only): F1 = 97.62 ± 0.00
- NuExtract:2B (NuExtract, Text-only): F1 = 97.62 ± 0.00

Bottom 5 Performers:

- granite3.32b (Granite, Text-only): F1 = 37.28 ± 41.09
- qwen31.7b (Qwen, Text-only): F1 = 43.51 ± 26.89
- llama3.21b (Llama, Text-only): F1 = 46.06 ± 21.61
- gpt-4o (GPT, with Vision): F1 = 49.58 ± 23.24
- granite3.2 (Granite, with Vision): F1 = 49.80 ± 18.49

Error Analysis Summary:

- Average False Positives: 14.80
- Average False Negatives: 26.49
- Models with more FP than FN: 20
- Models with more FN than FP: 109

SOURCE CATEGORY DETAILED STATISTICS:

COMMERCIAL MODELS:

GPT:

- Accuracy: 47.2±14.8 (n=46)
- F1 Score: 59.2±15.7
- Precision: 75.6±21.1
- False Positives: 9.5
- False Negatives: 25.6

NuExtract:

- Accuracy: 97.6±0.0 (n=8)
- F1 Score: 97.6±0.0
- Precision: 97.6±0.0
- False Positives: 1.0
- False Negatives: 1.0

OLLAMA MODELS:

Gemma/Gemini:

- Accuracy: 48.5±11.0 (n=29)
- F1 Score: 58.1±11.3
- Precision: 63.9±15.4
- False Positives: 16.8
- False Negatives: 24.3

Granite:

- Accuracy: 35.2±22.9 (n=19)
- F1 Score: 47.8±22.2
- Precision: 60.6±37.0
- False Positives: 17.2
- False Negatives: 31.9

Llama:

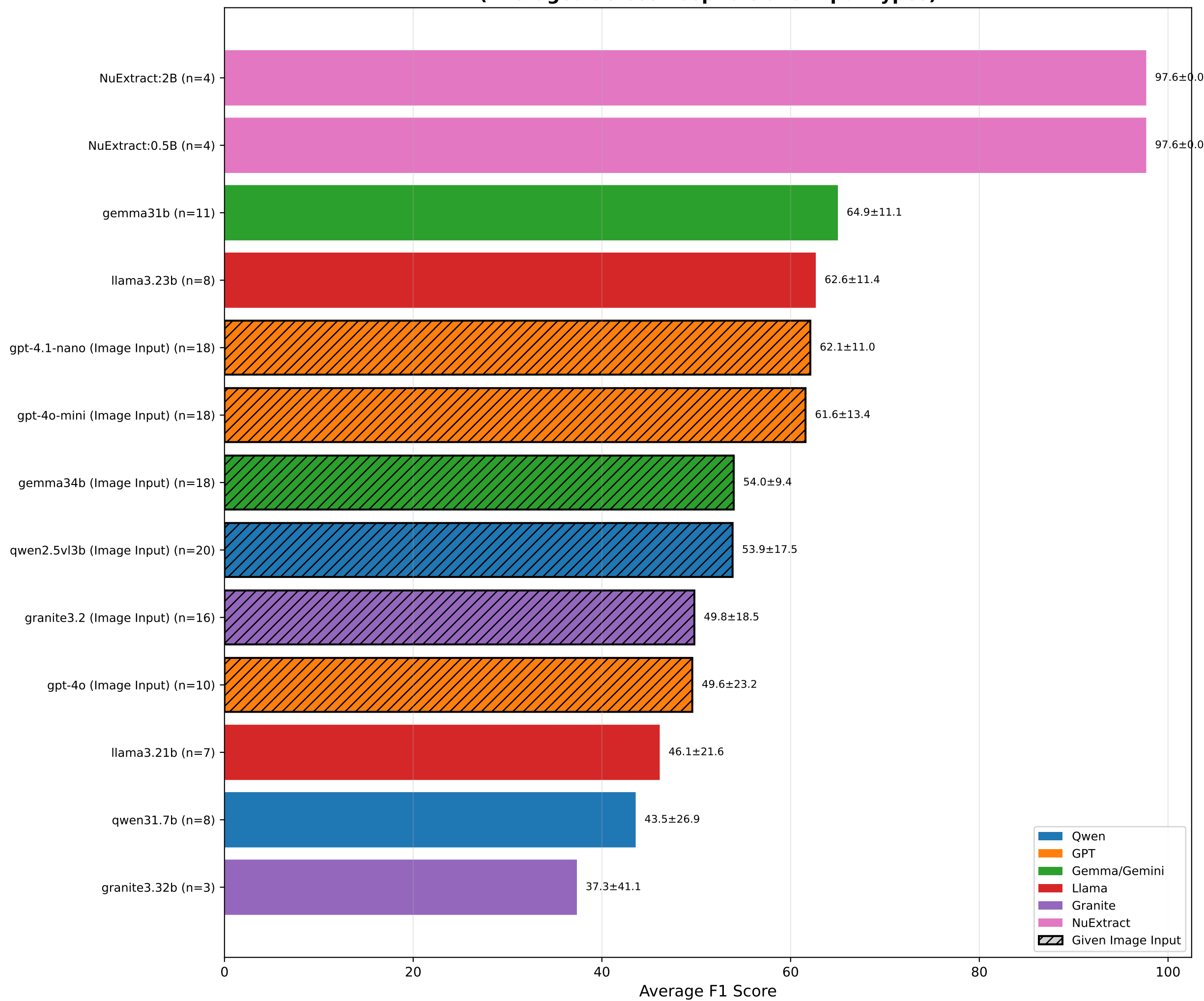
- Accuracy: 42.7±15.9 (n=15)
- F1 Score: 54.9±18.4
- Precision: 69.9±22.6
- False Positives: 15.0
- False Negatives: 31.7

Qwen:

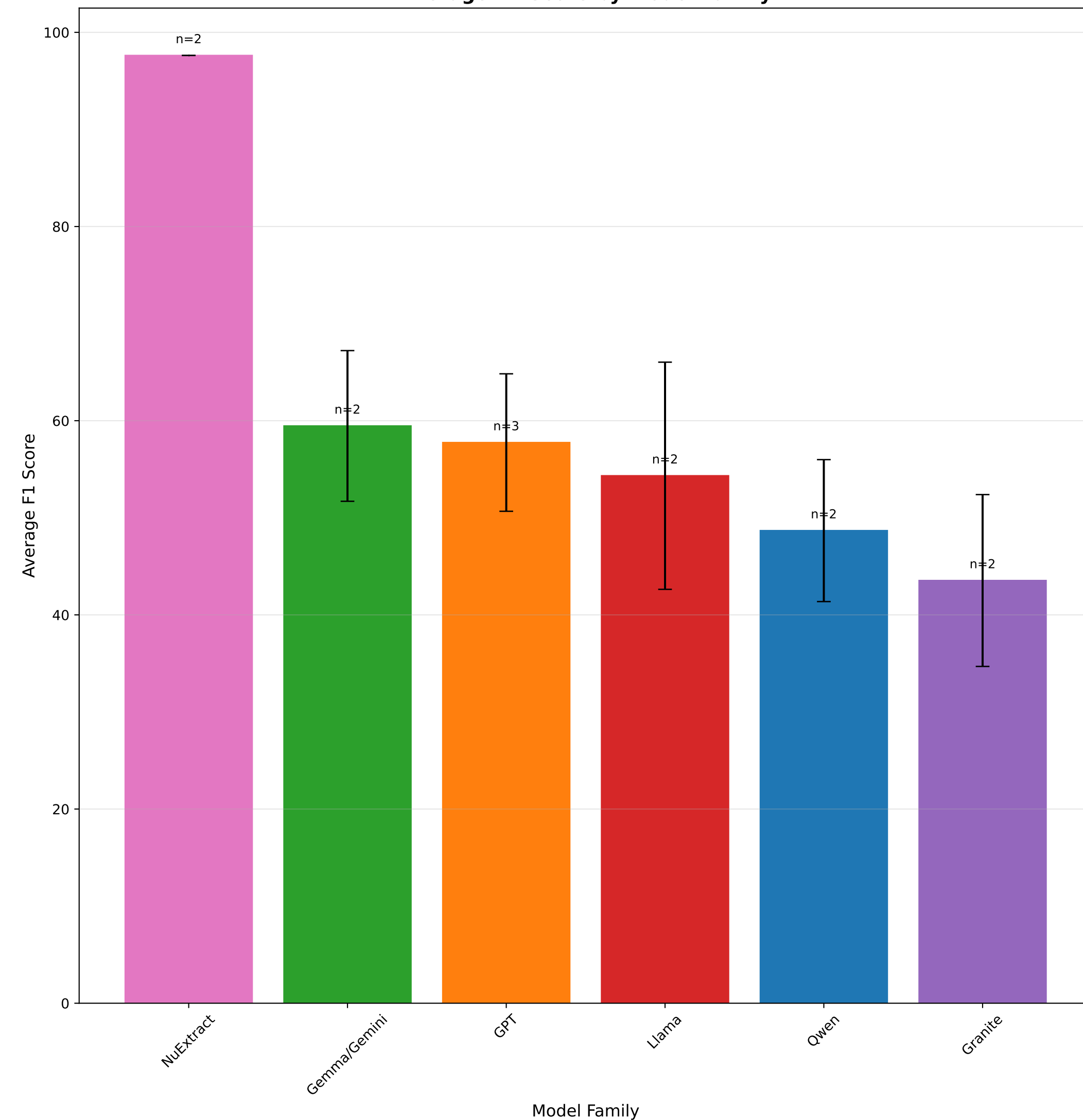
- Accuracy: 42.3±19.1 (n=28)
- F1 Score: 50.9±20.6
- Precision: 55.8±22.5
- False Positives: 23.8
- False Negatives: 31.0

=====

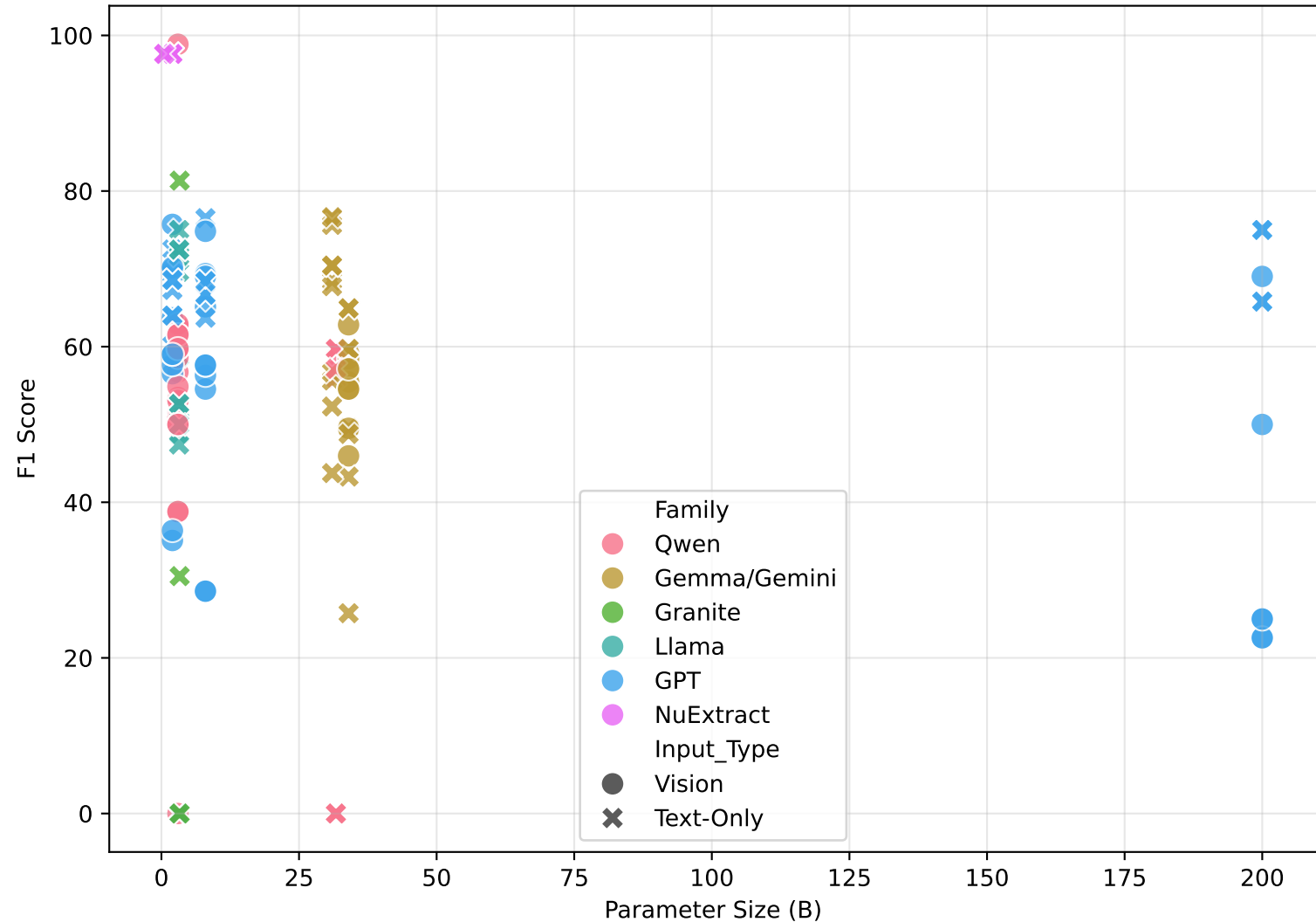
Overall F1 Score Performance - Models Grouped by Base Name  
(Averaged across hospitals and input types)



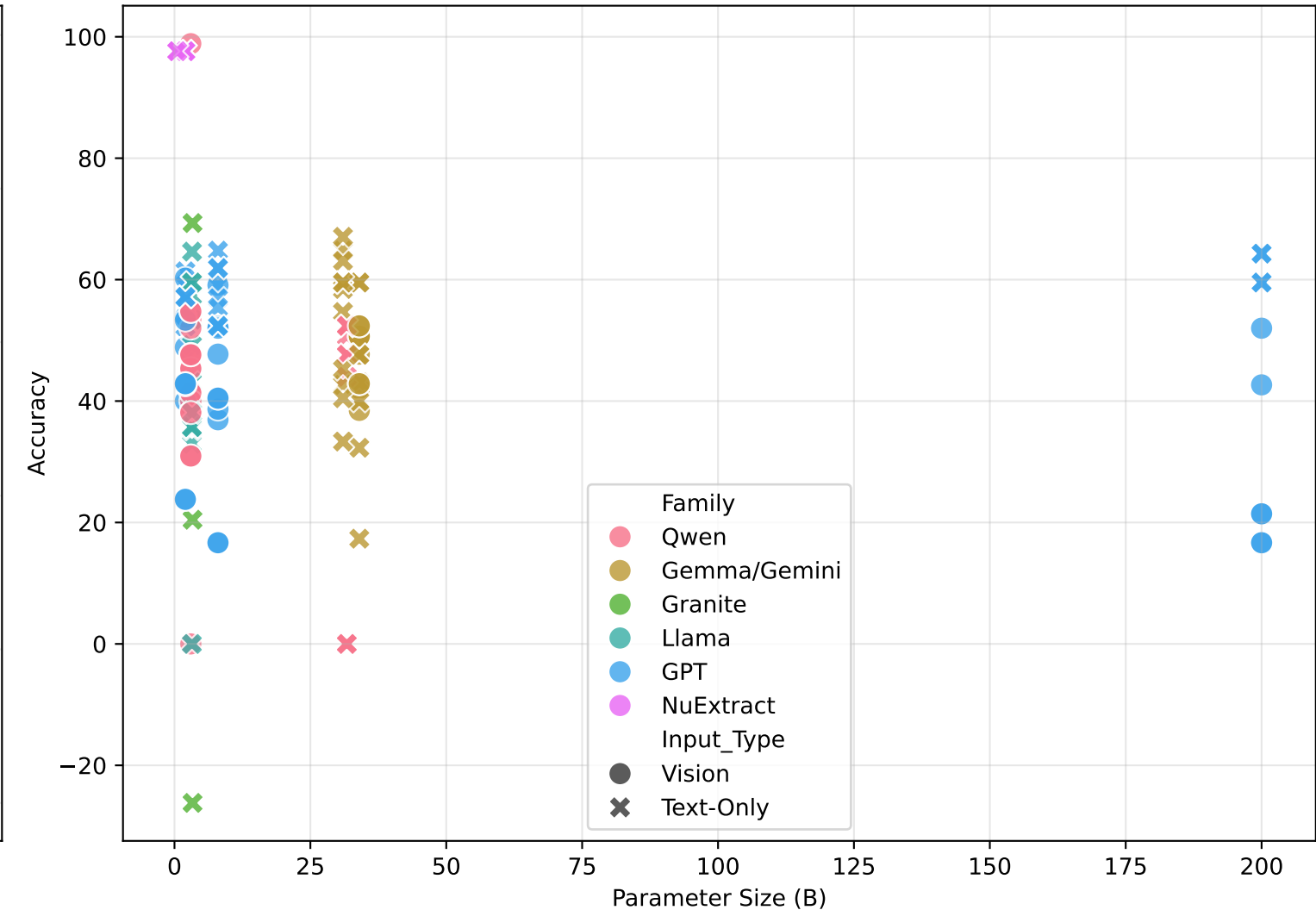
Average F1 Score by Model Family



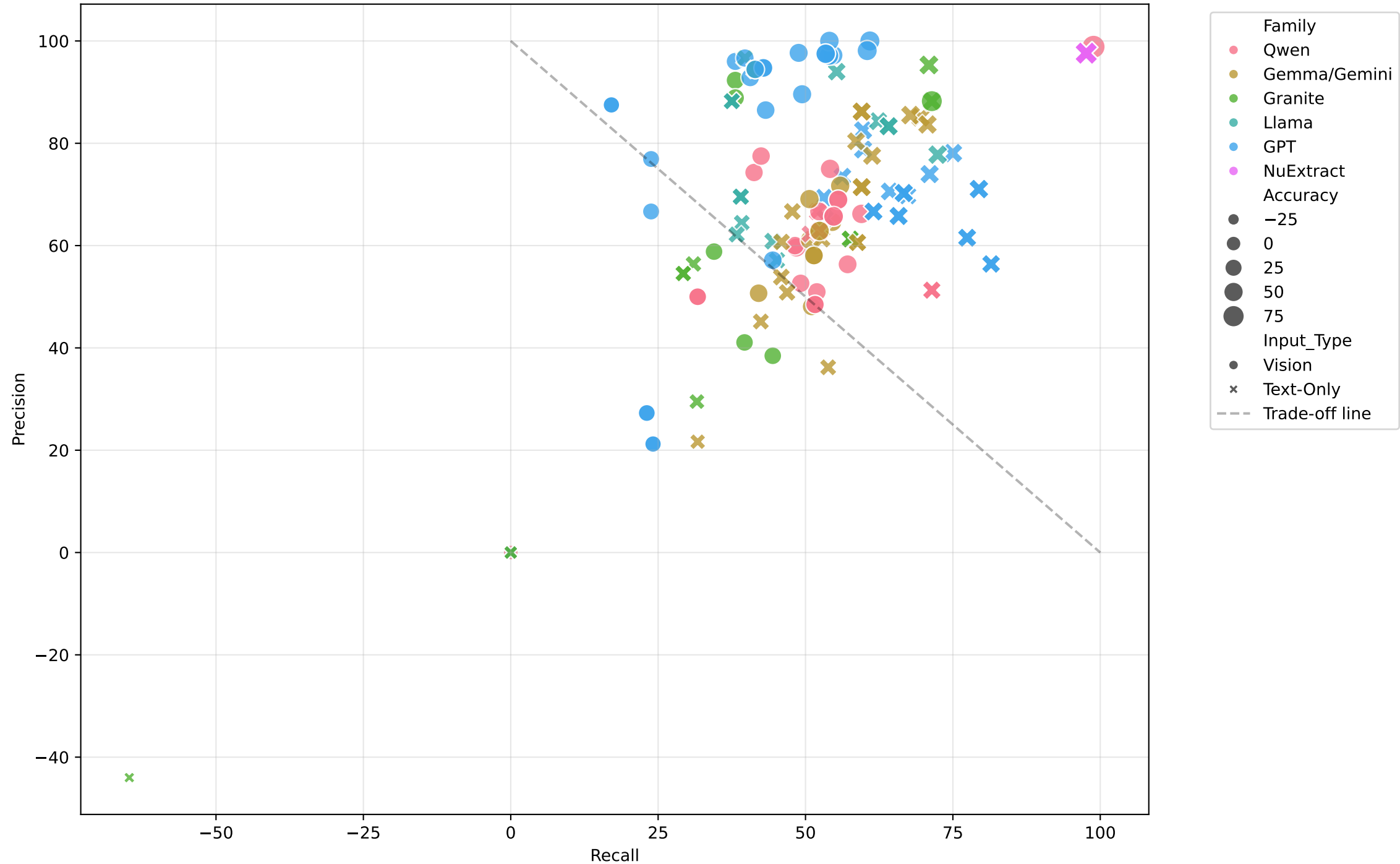
### F1 Score vs Parameter Size by Family and Input Type

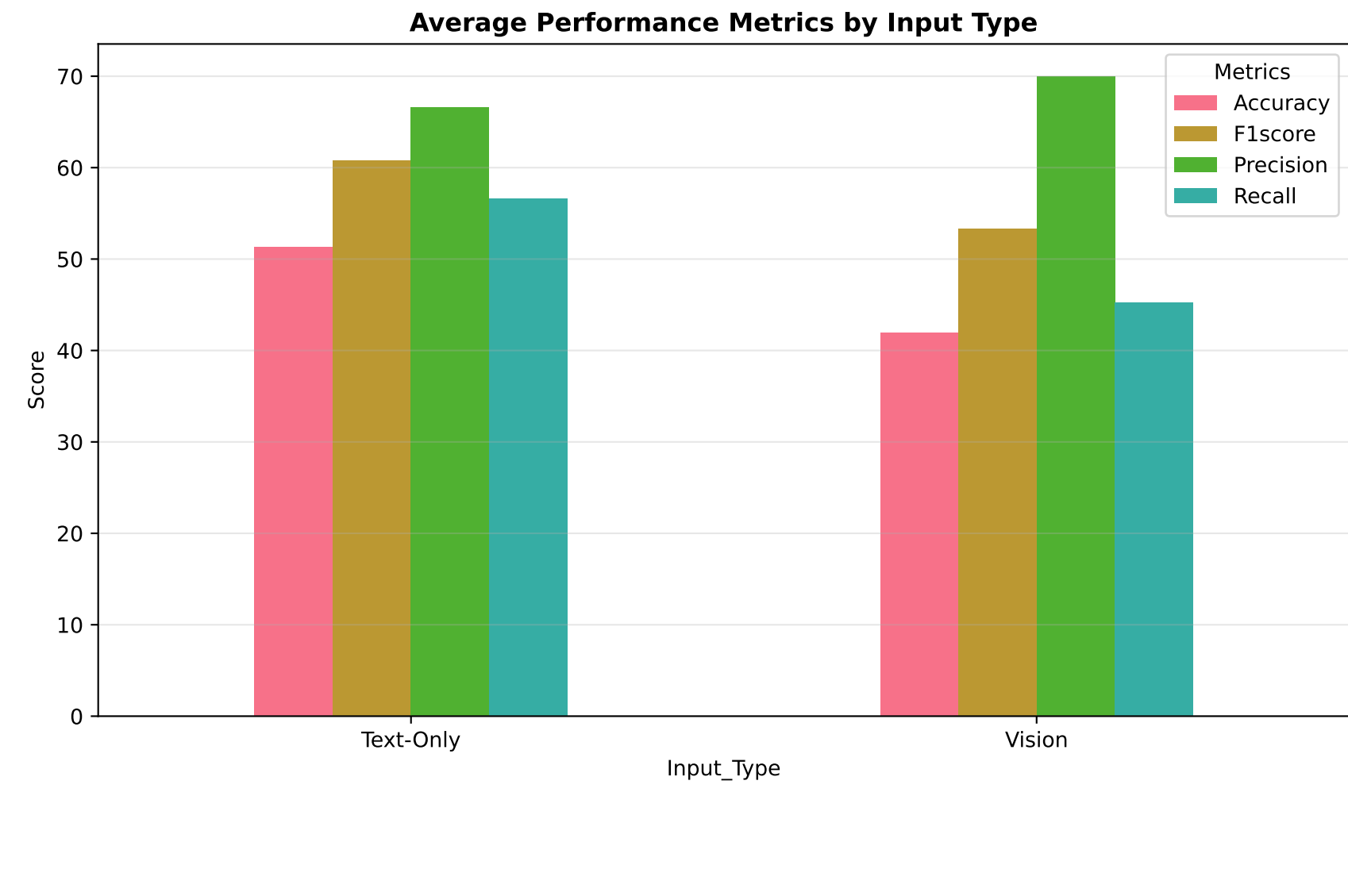
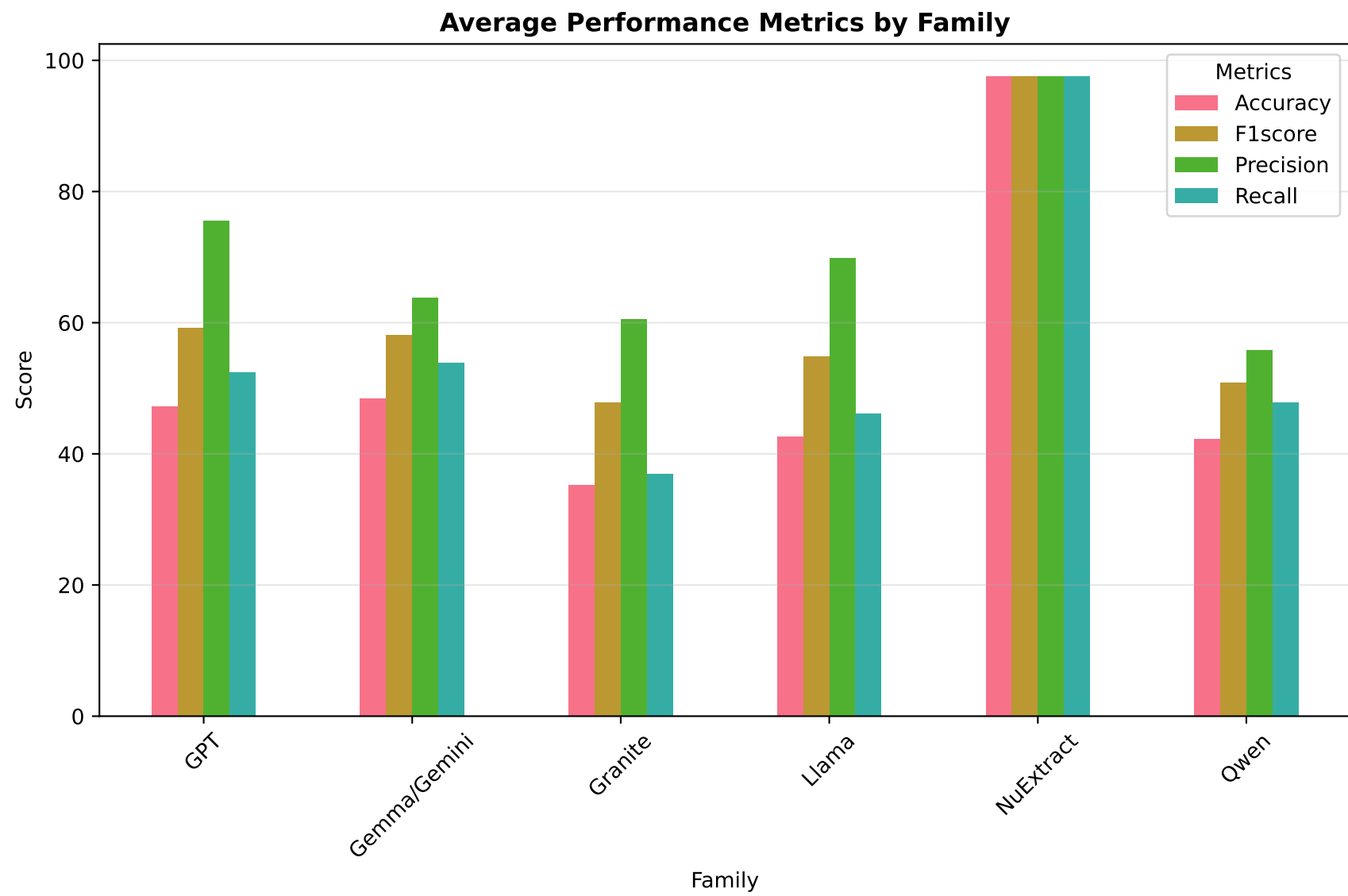
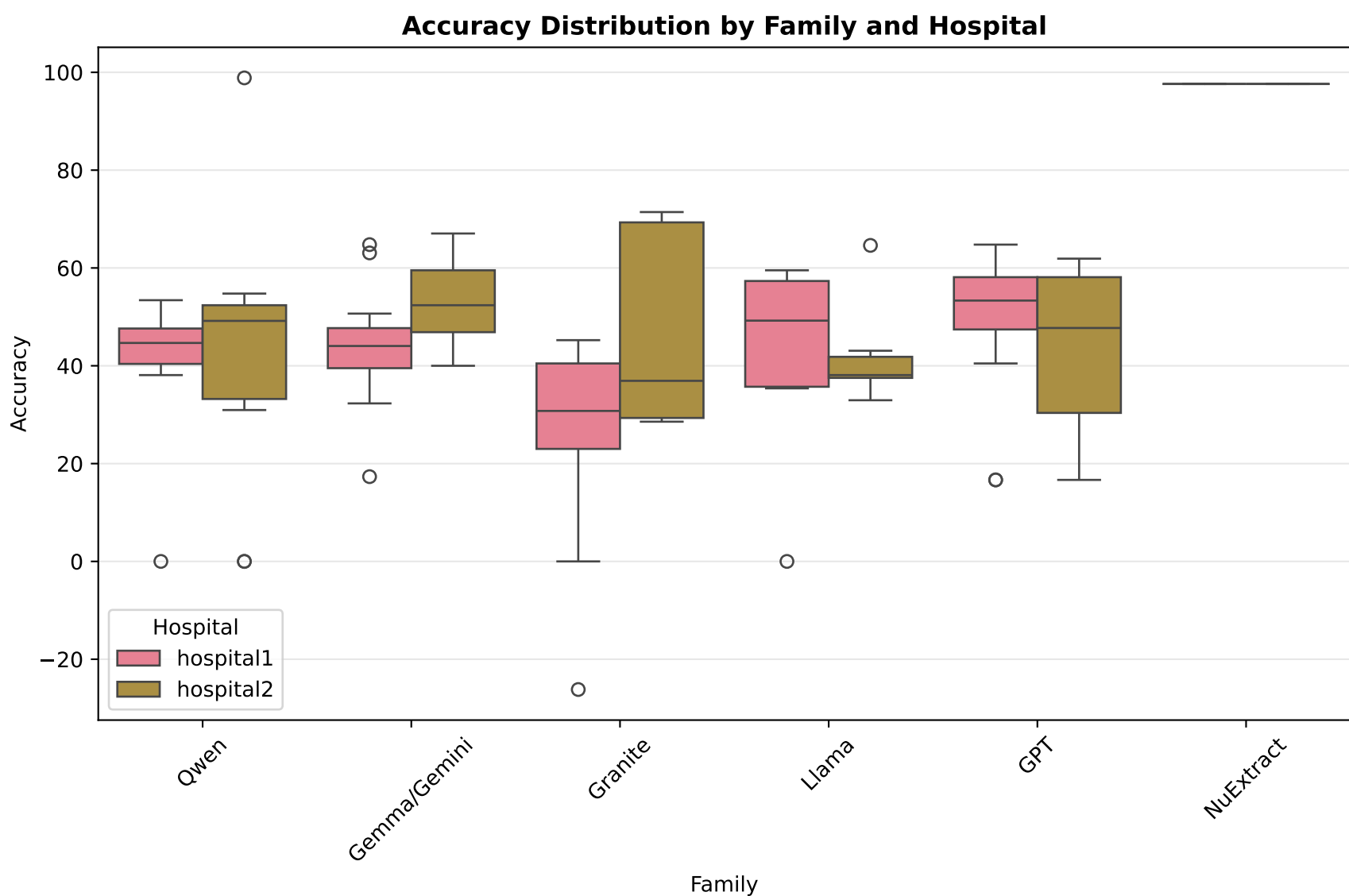
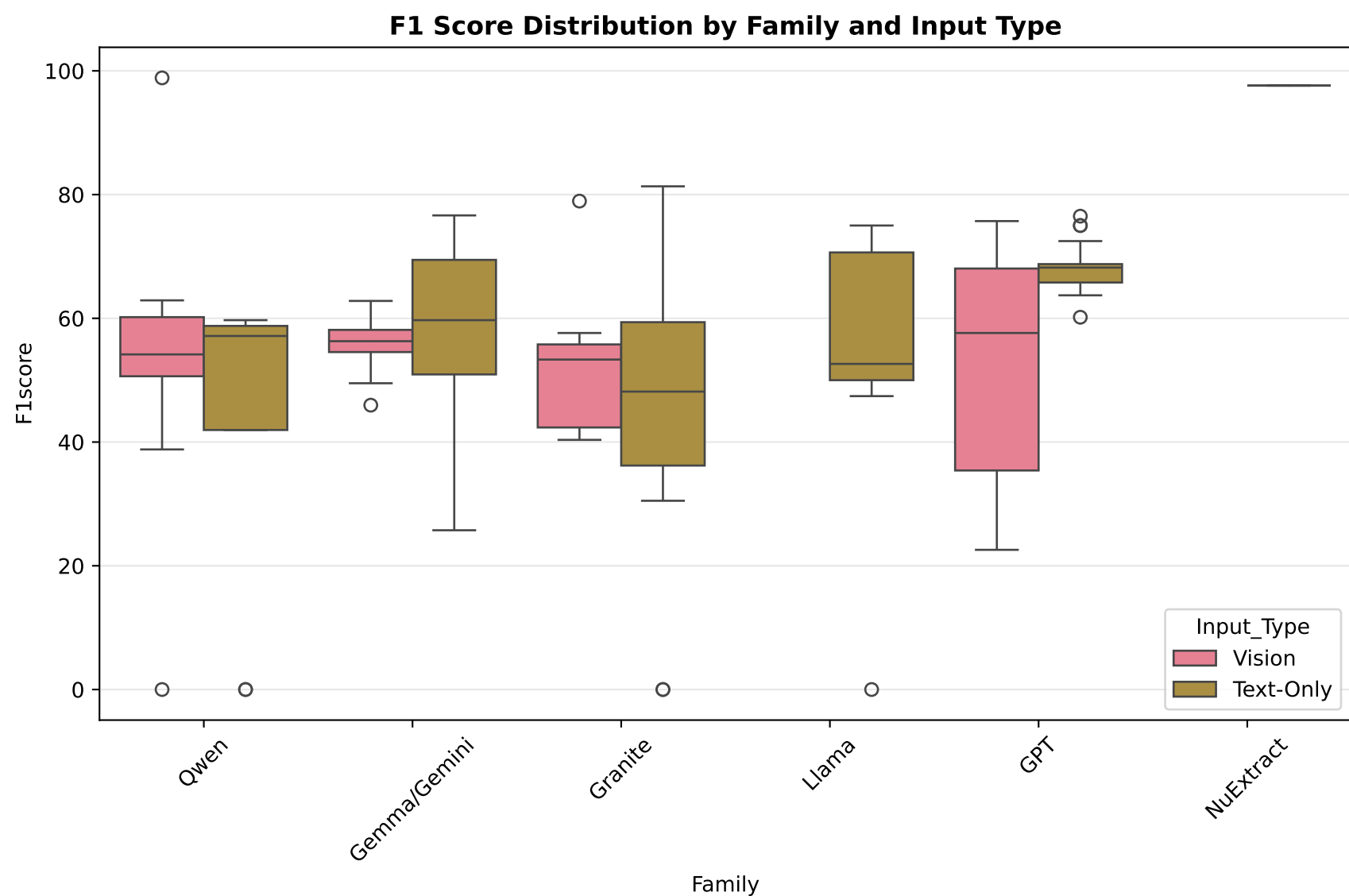


### Accuracy vs Parameter Size by Family and Input Type

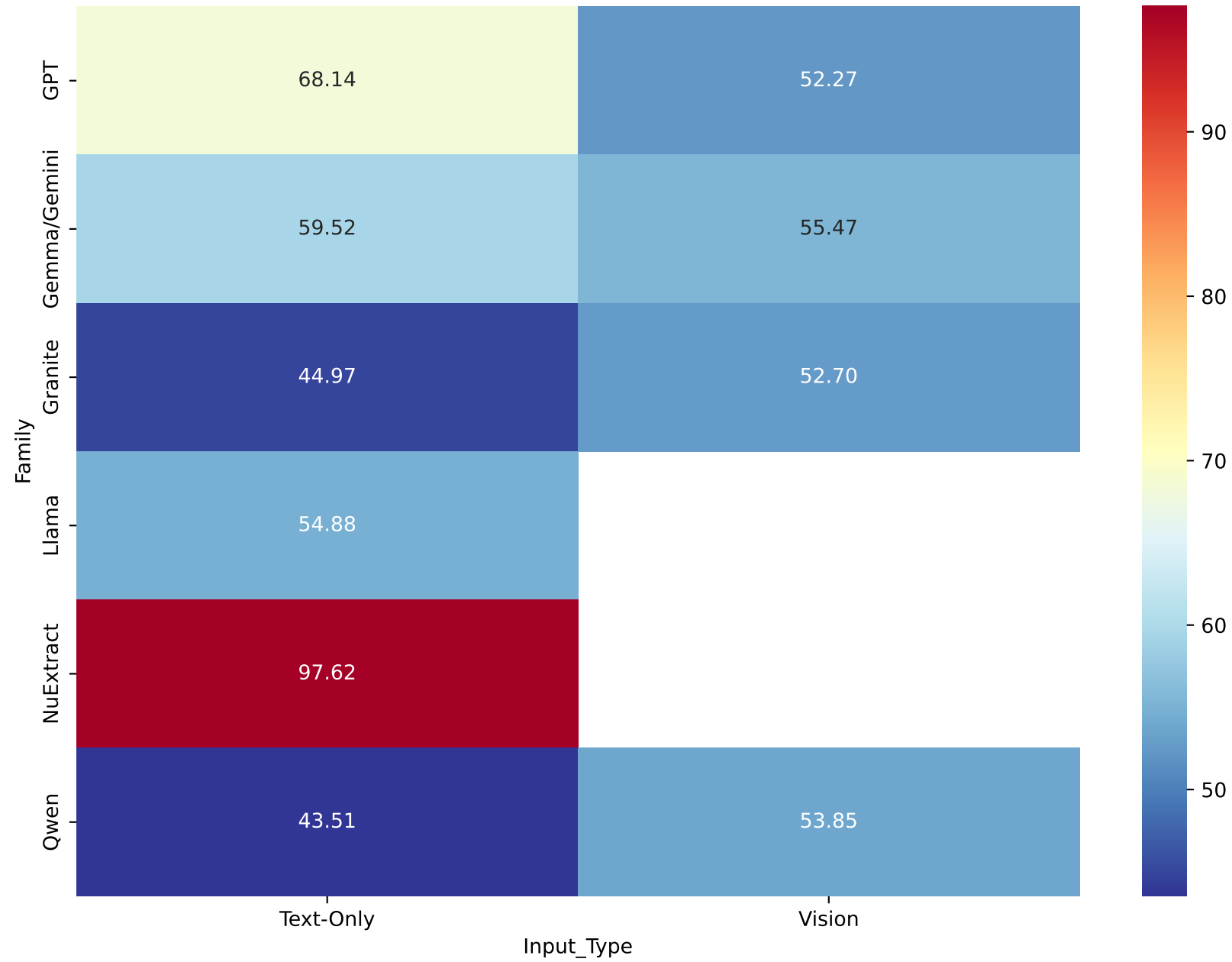


**Precision vs Recall by Family and Input Type**  
(Size = Accuracy, Color = Family)

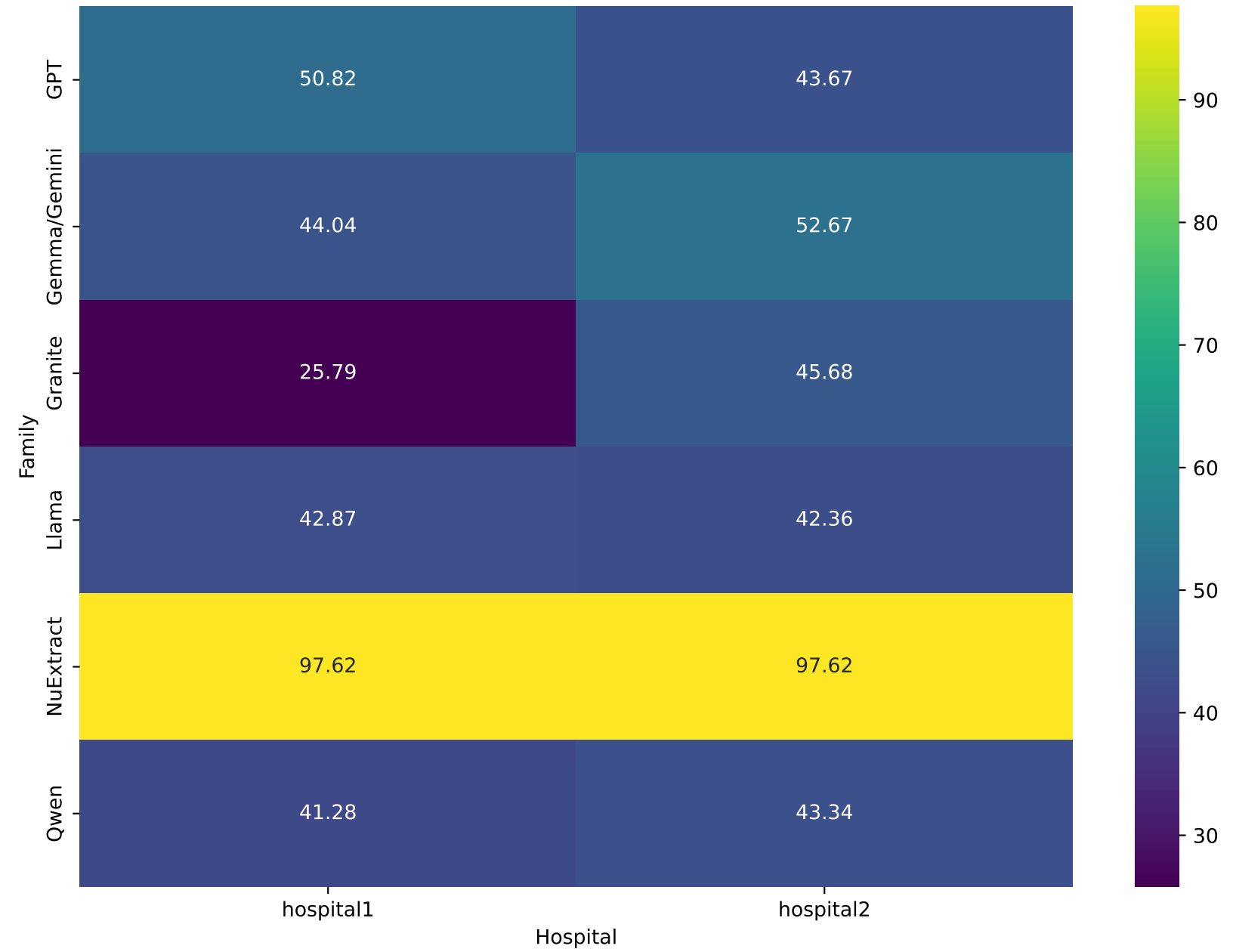




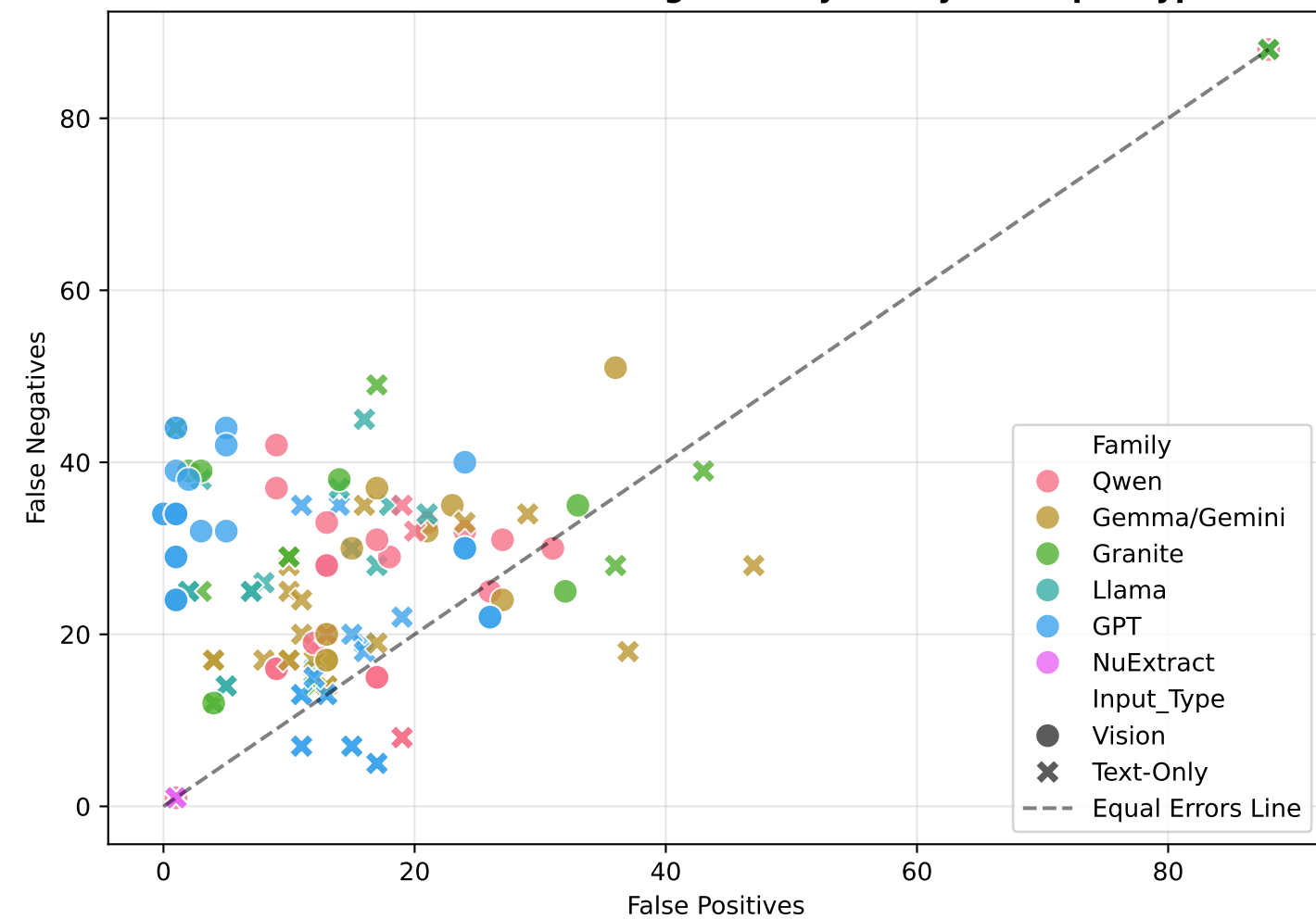
Average F1 Score: Family vs Input Type



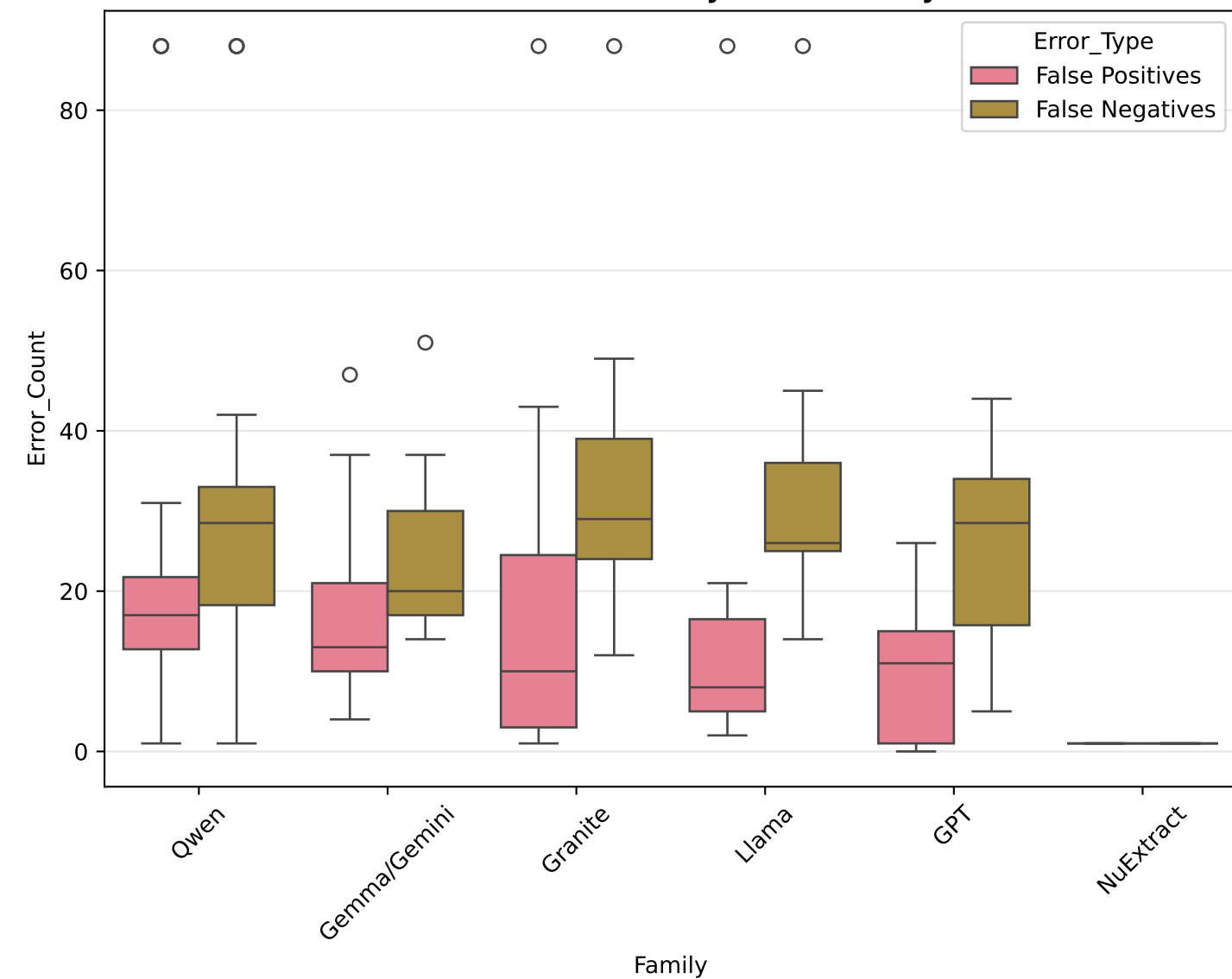
Average Accuracy: Family vs Hospital



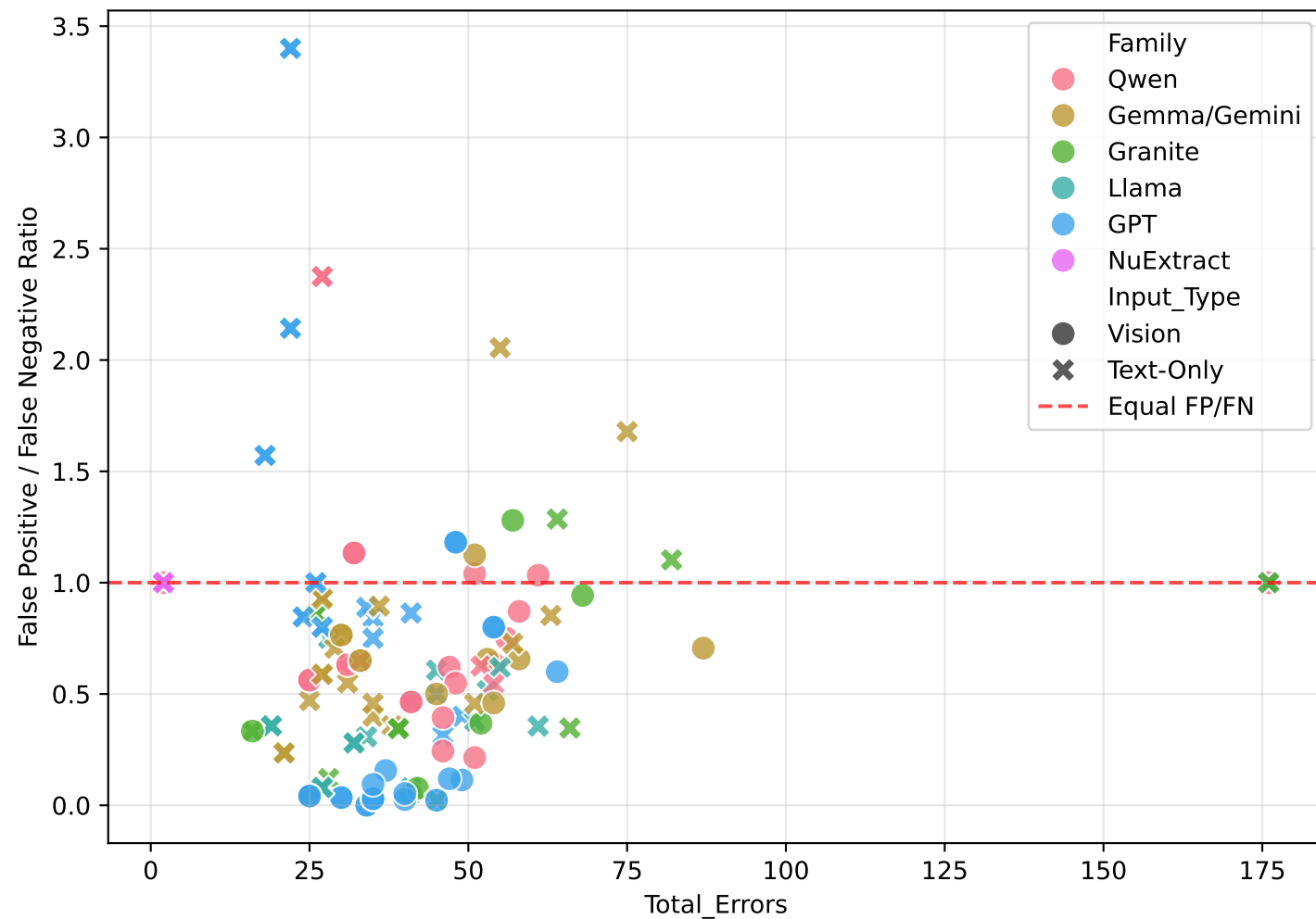
False Positives vs False Negatives by Family and Input Type



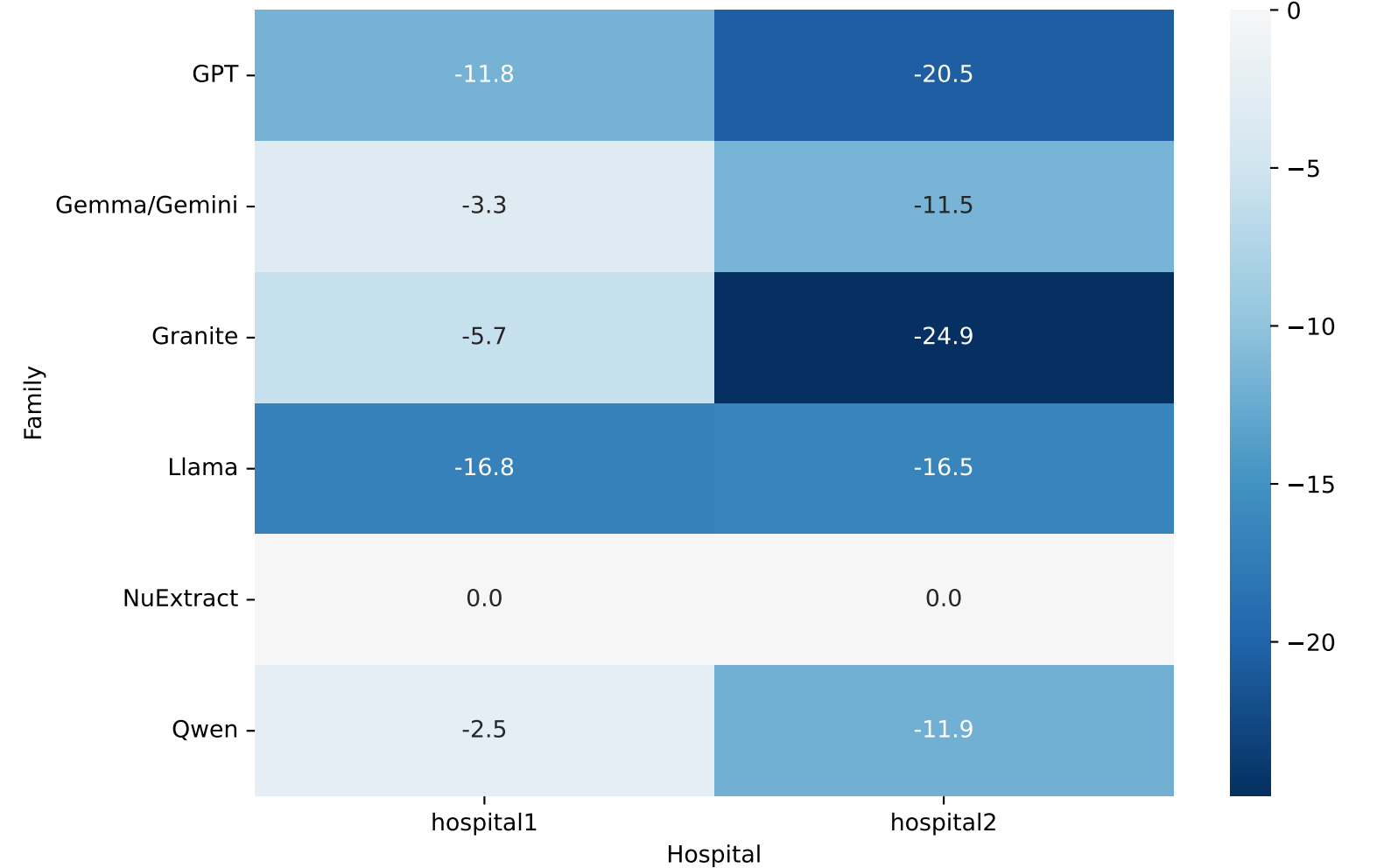
Error Distribution by Model Family

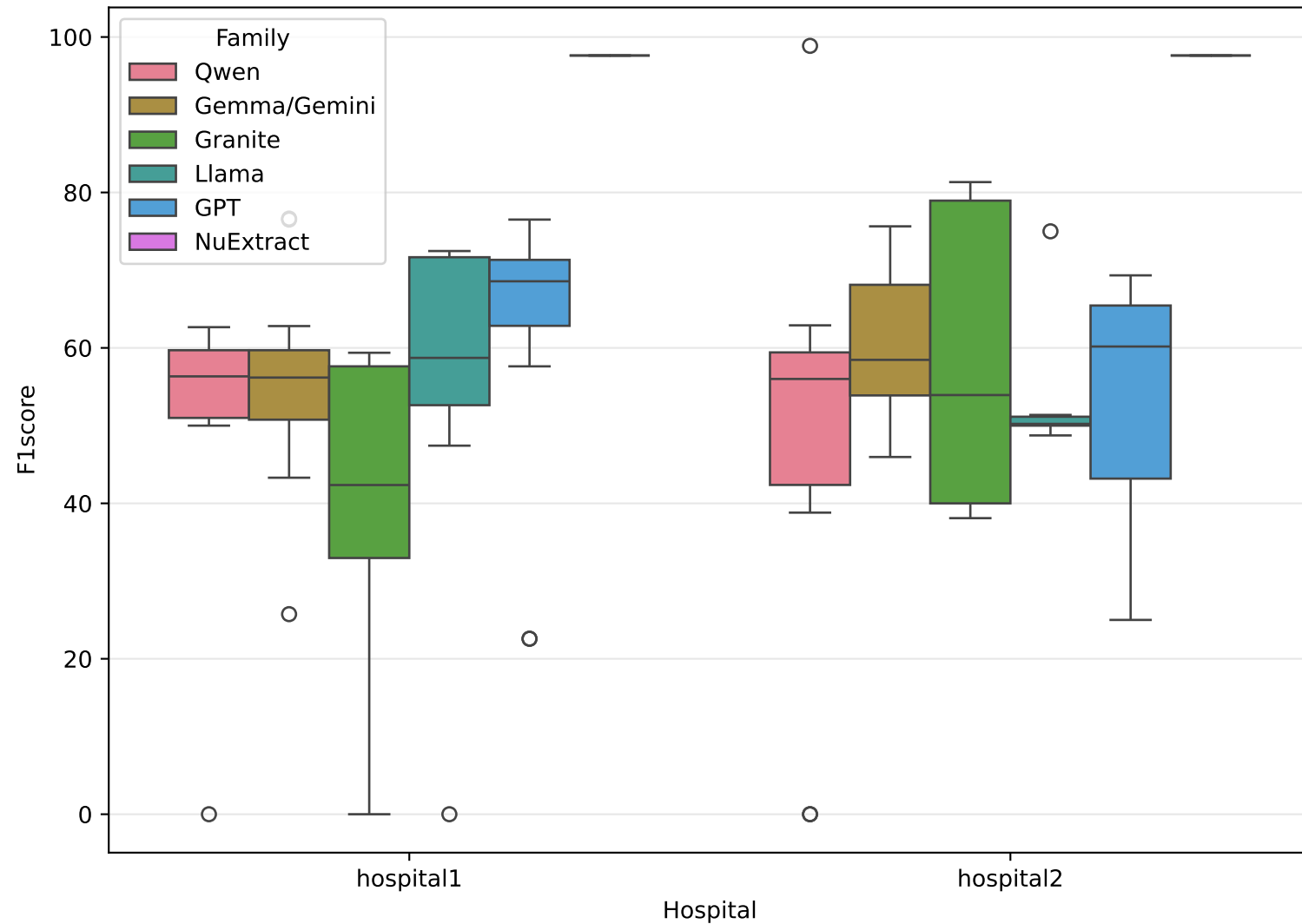
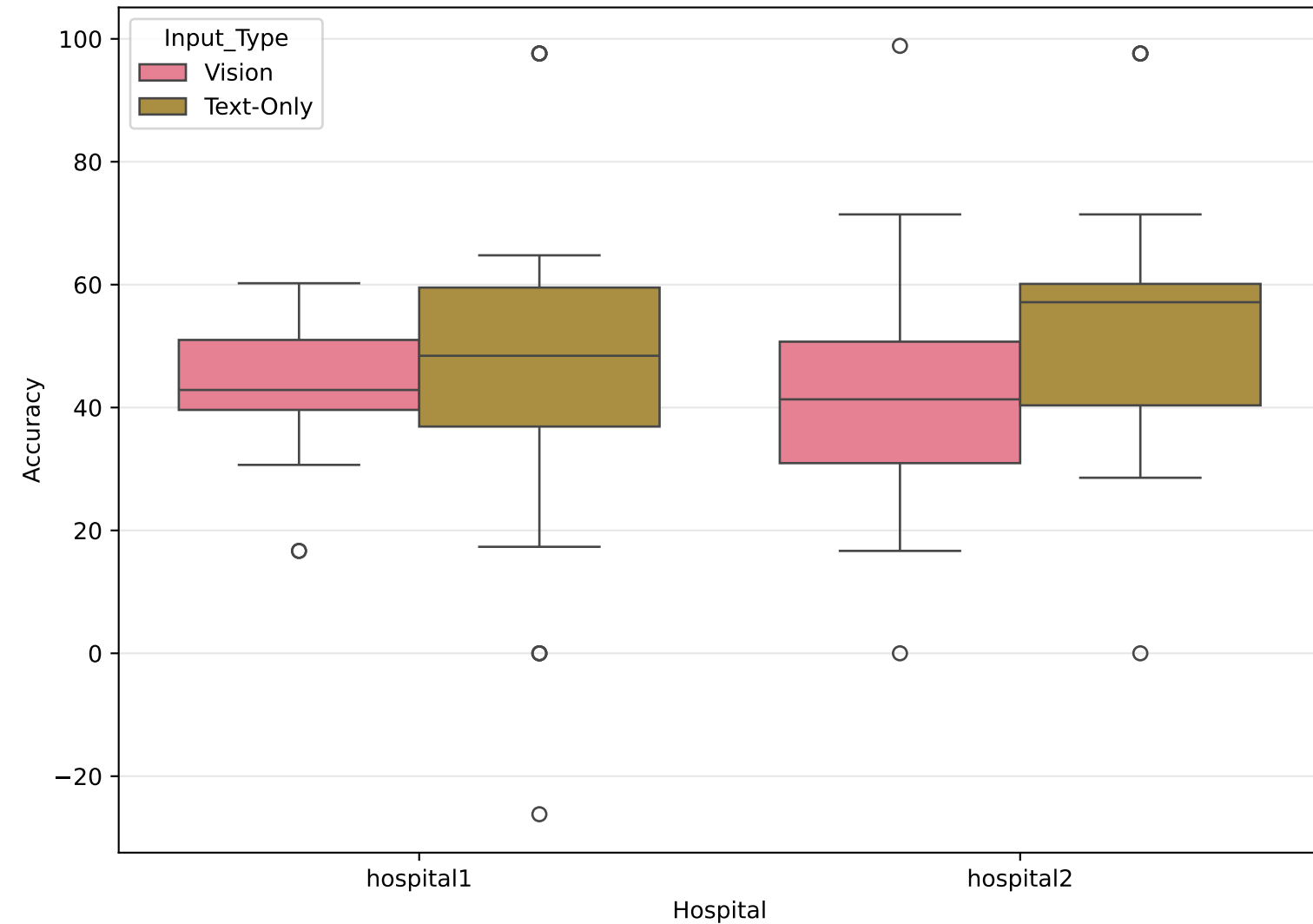


Error Ratio vs Total Errors  
(Ratio = FP/FN)



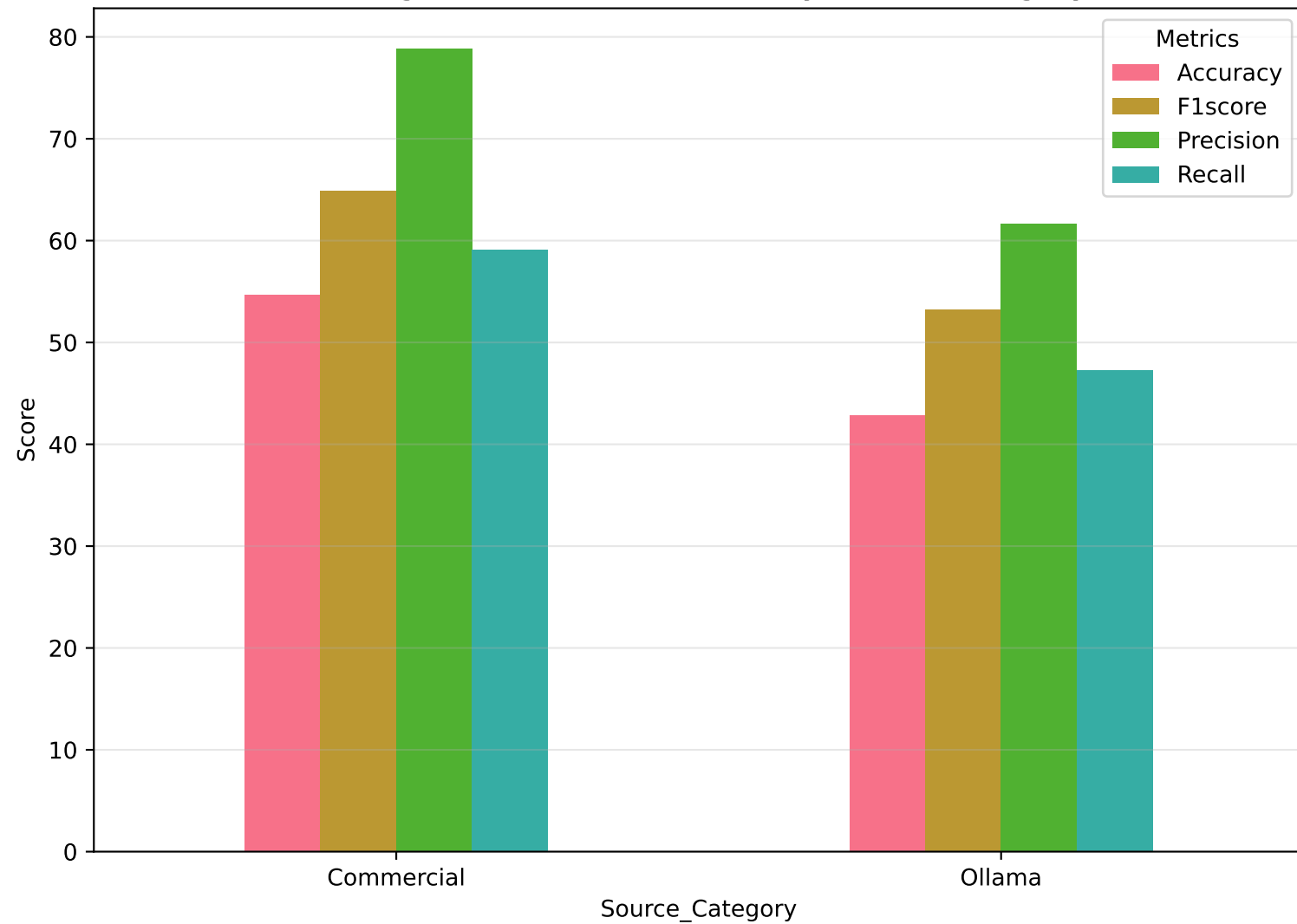
Error Bias: FP - FN by Family and Hospital  
(+ve = More FP, -ve = More FN)



**F1 Score Distribution by Hospital and Family****Accuracy Distribution by Hospital and Input Type**



**Average Performance Metrics by Source Category**



**F1 Score Distribution by Source and Input Type**

