```
COMPREHENSIVE LLM PERFORMANCE ANALYSIS
=======================================
Report Generated: 2025-06-26 15:32:30

OVERALL PERFORMANCE METRICS:
    • Average F1 Score: 58.082 (Std: 18.227)
    • Average Accuracy: 43.149 (Std: 16.620)
    • Average Precision: 61.895 (Std: 21.505)
    • Average Recall: 57.942 (Std: 18.549)

GROUPED MODEL F1 SCORE STATISTICS:
    • Unique Base Models: 13
    • Total Test Instances: 76
    • Best Performing Model: NuExtract:2B (F1: 75.16)
    • Worst Performing Model: GLiNER:NuNerZero (F1: 18.70)
    • Overall Average F1: 57.83
    • Models with Vision: 6
```
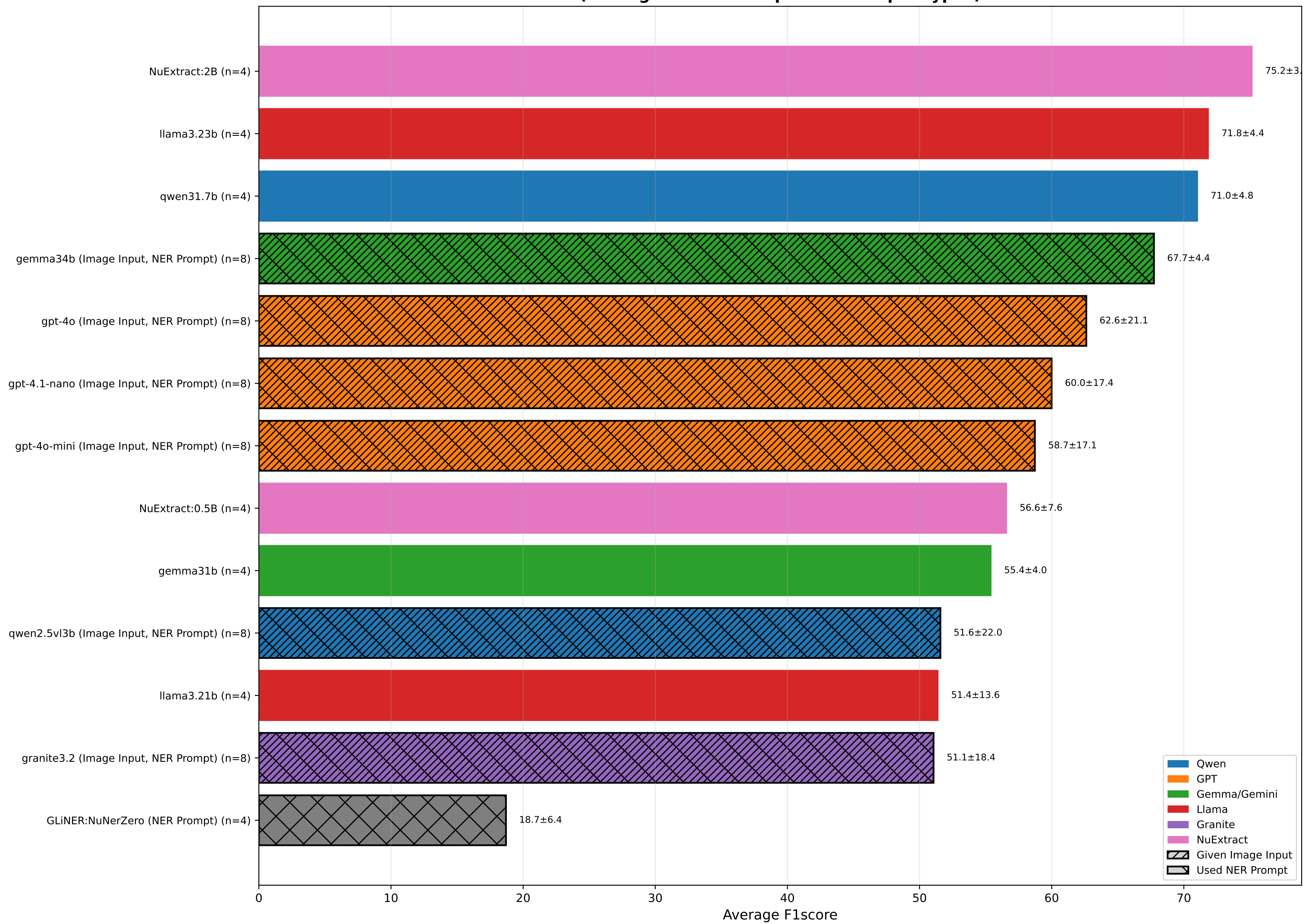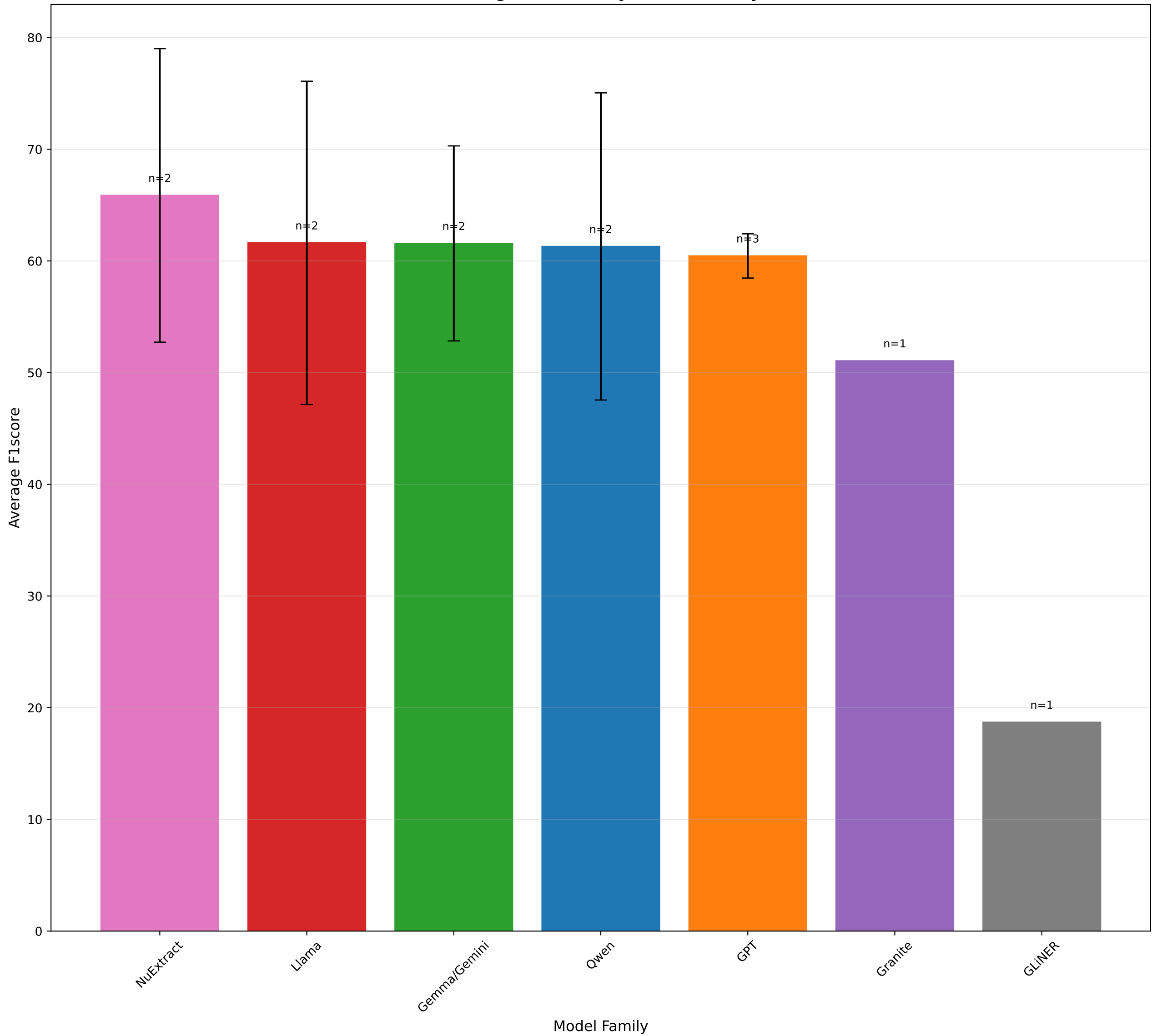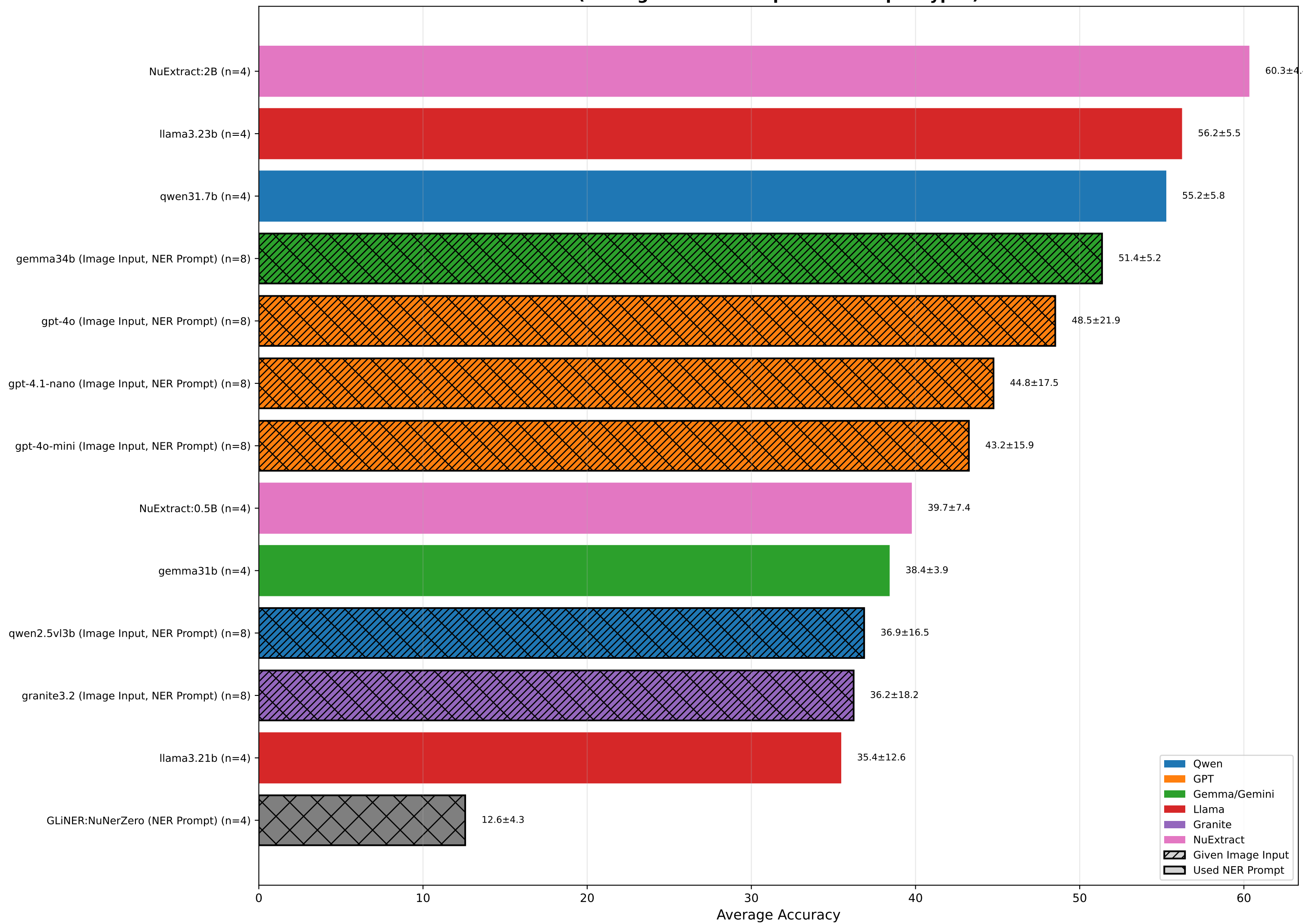
**Overall F1score Performance - Models Grouped by Base Name**
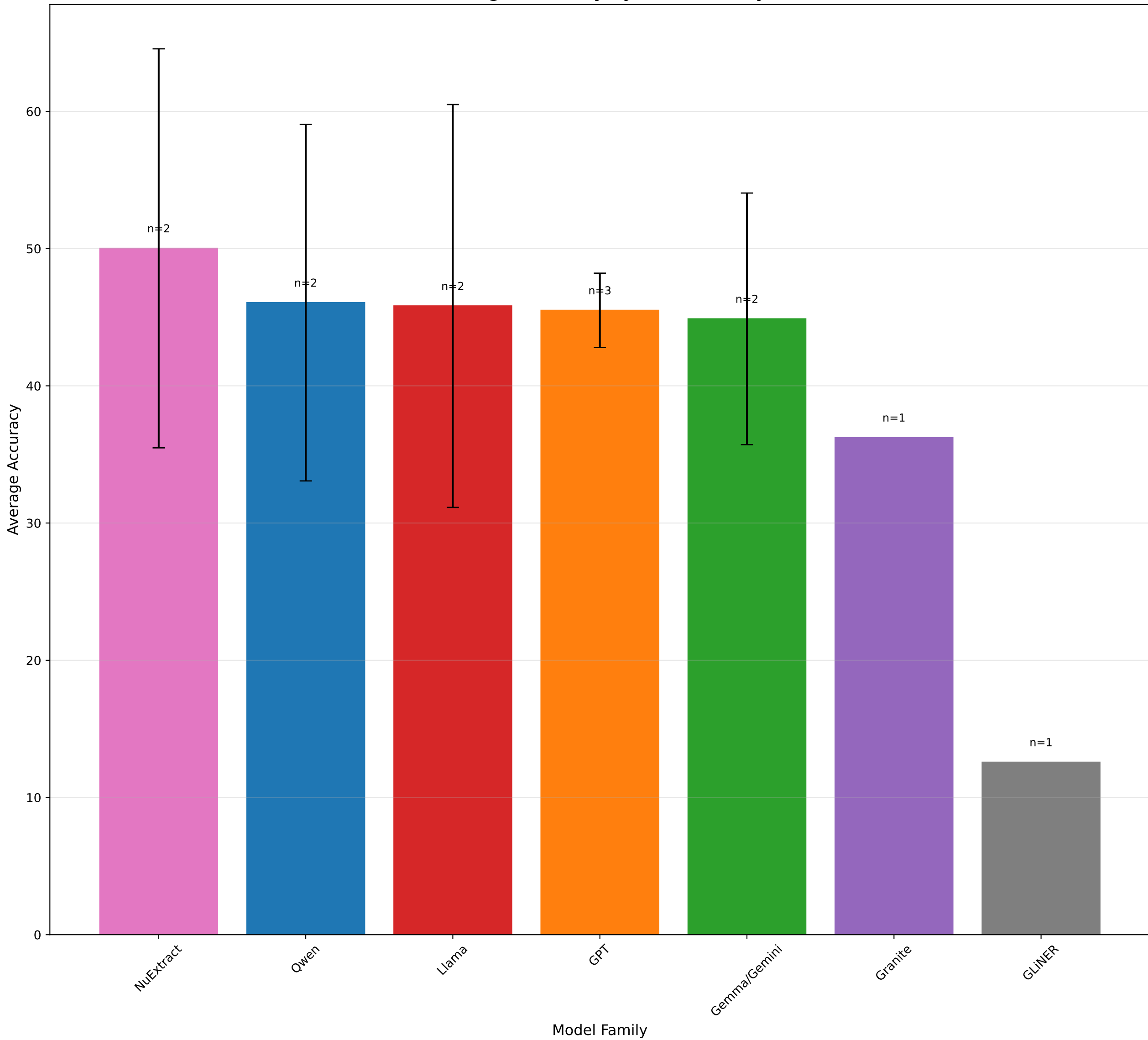**(Averaged across hospitals and input types)**

| Model | Average F1score |
|---|---|
| NuExtract:2B (n=4) | 75.2±3.5 |
| llama3.23b (n=4) | 71.8±4.4 |
| qwen31.7b (n=4) | 71.0±4.8 |
| gemma34b (Image Input, NER Prompt) (n=8) | 67.7±4.4 |
| gpt-4o (Image Input, NER Prompt) (n=8) | 62.6±21.1 |
| gpt-4.1-nano (Image Input, NER Prompt) (n=8) | 60.0±17.4 |
| gpt-4o-mini (Image Input, NER Prompt) (n=8) | 58.7±17.1 |
| NuExtract:0.5B (n=4) | 56.6±7.6 |
| gemma31b (n=4) | 55.4±4.0 |
| qwen2.5vl3b (Image Input, NER Prompt) (n=8) | 51.6±22.0 |
| llama3.21b (n=4) | 51.4±13.6 |
| granite3.2 (Image Input, NER Prompt) (n=8) | 51.1±18.4 |
| GLiNER:NuNerZero (NER Prompt) (n=4) | 18.7±6.4 |

Legend: Qwen, GPT, Gemma/Gemini, Llama, Granite, NuExtract, Given Image Input, Used NER Prompt

**Average F1score by Model Family**

| Model Family | Average F1score |
|---|---|
| NuExtract | n=2 |
| Llama | n=2 |
| Gemma/Gemini | n=2 |
| Qwen | n=2 |
| GPT | n=3 |
| Granite | n=1 |
| GLINER | n=1 |

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

| Model | Average Accuracy |
|---|---|
| NuExtract:2B (n=4) | 60.3±4.4 |
| llama3.23b (n=4) | 56.2±5.5 |
| qwen31.7b (n=4) | 55.2±5.8 |
| gemma34b (Image Input, NER Prompt) (n=8) | 51.4±5.2 |
| gpt-4o (Image Input, NER Prompt) (n=8) | 48.5±21.9 |
| gpt-4.1-nano (Image Input, NER Prompt) (n=8) | 44.8±17.5 |
| gpt-4o-mini (Image Input, NER Prompt) (n=8) | 43.2±15.9 |
| NuExtract:0.5B (n=4) | 39.7±7.4 |
| gemma31b (n=4) | 38.4±3.9 |
| qwen2.5vl3b (Image Input, NER Prompt) (n=8) | 36.9±16.5 |
| granite3.2 (Image Input, NER Prompt) (n=8) | 36.2±18.2 |
| llama3.21b (n=4) | 35.4±12.6 |
| GLiNER:NuNerZero (NER Prompt) (n=4) | 12.6±4.3 |

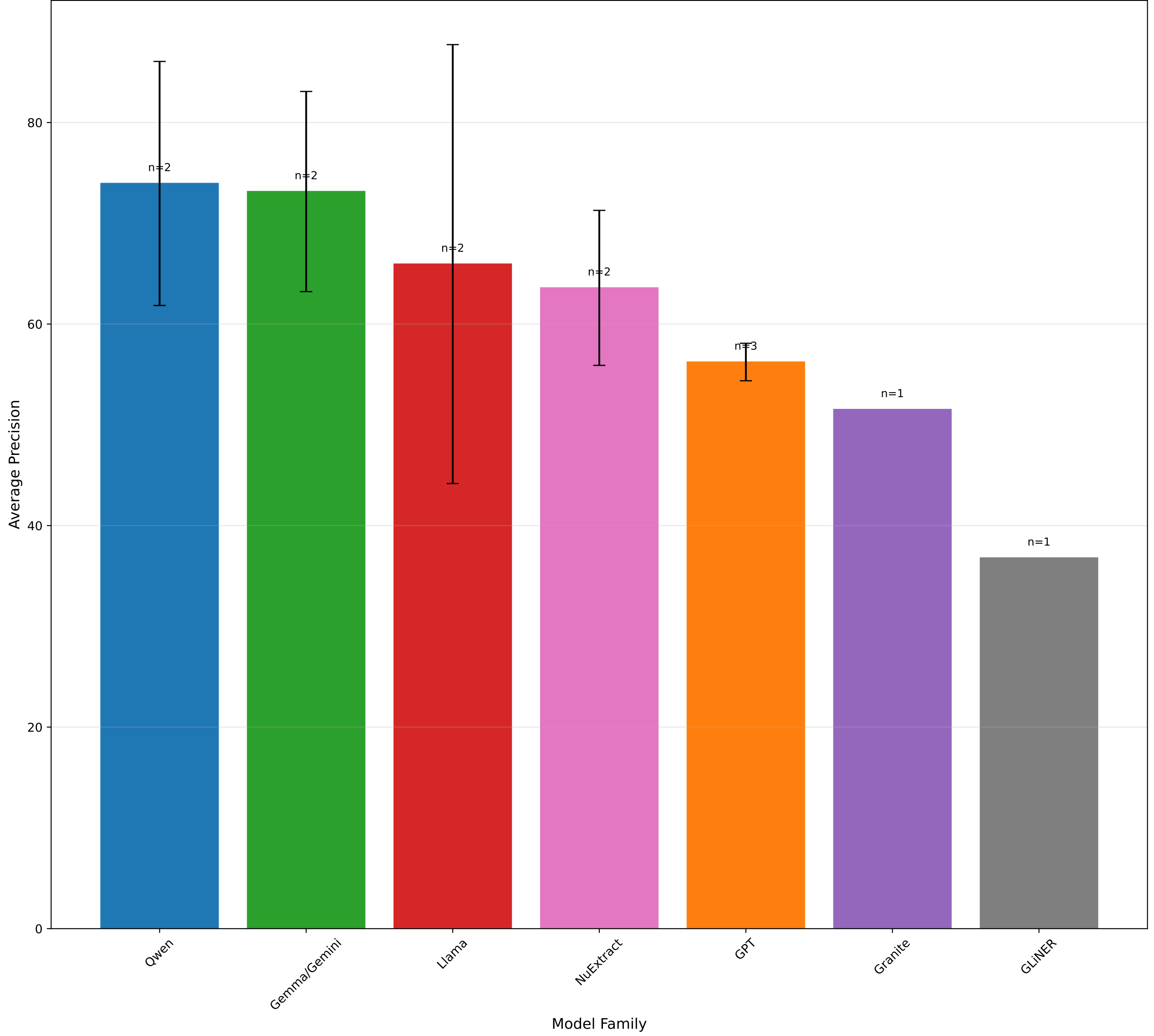Legend: Qwen, GPT, Gemma/Gemini, Llama, Granite, NuExtract, Given Image Input, Used NER Prompt

**Average Accuracy by Model Family**

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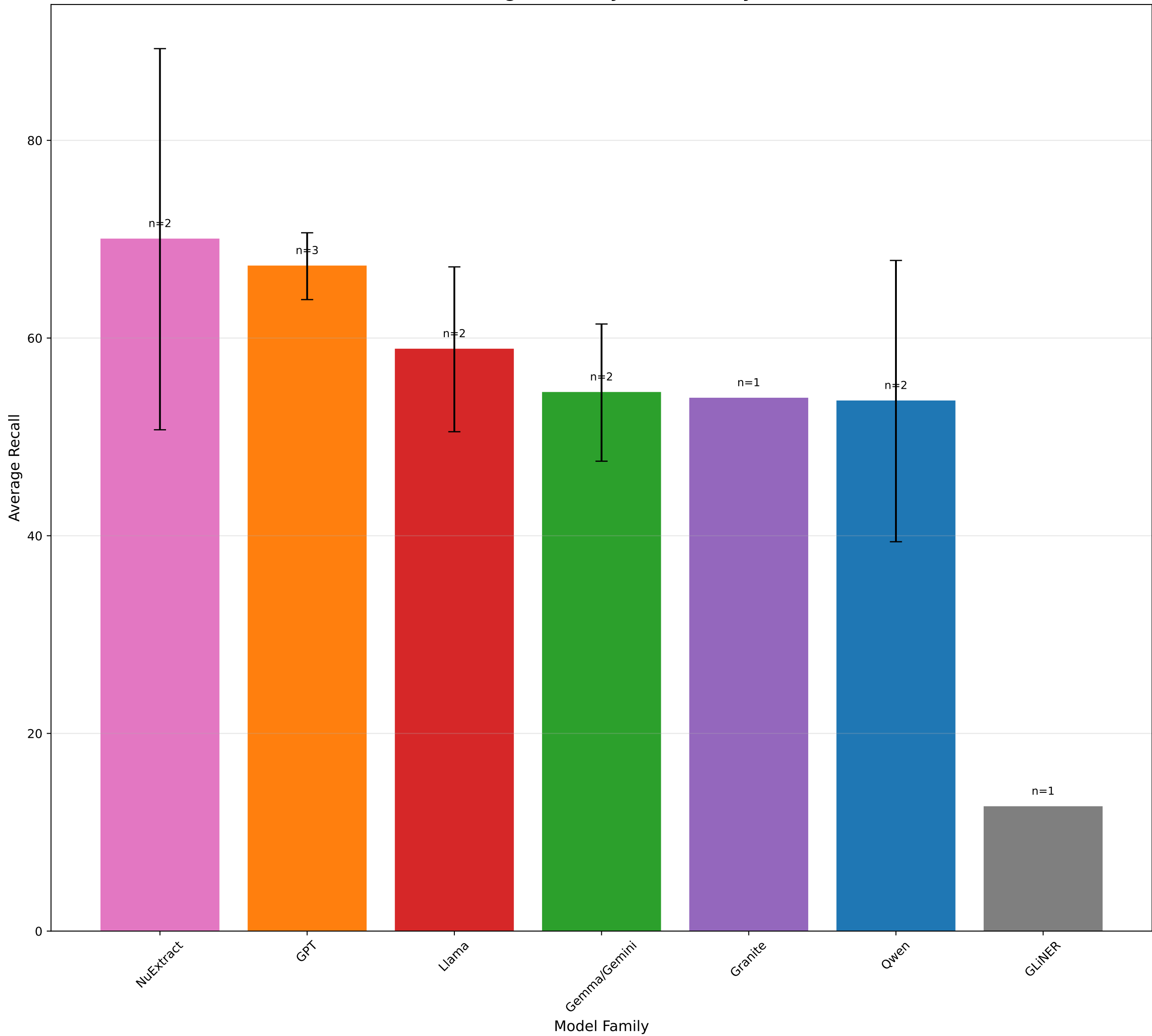**(Averaged across hospitals and input types)**

- qwen31.7b (n=4): 82.5±15.8
- llama3.23b (n=4): 81.3±4.6
- gemma34b (Image Input, NER Prompt) (n=8): 80.2±11.5
- NuExtract:2B (n=4): 69.0±8.4
- gemma31b (n=4): 66.1±10.0
- qwen2.5vl3b (Image Input, NER Prompt) (n=8): 65.4±29.8
- NuExtract:0.5B (n=4): 58.1±10.4
- gpt-4.1-nano (Image Input, NER Prompt) (n=8): 57.7±20.6
- gpt-4o (Image Input, NER Prompt) (n=8): 56.8±22.2
- gpt-4o-mini (Image Input, NER Prompt) (n=8): 54.1±19.8
- granite3.2 (Image Input, NER Prompt) (n=8): 51.5±25.6
- llama3.21b (n=4): 50.5±14.0
- GLiNER:NuNerZero (NER Prompt) (n=4): 36.8±12.9

Legend:
- Qwen
- GPT
- Gemma/Gemini
- Llama
- Granite
- NuExtract
- Given Image Input
- Used NER Prompt

Average Precision

**Average Precision by Model Family**

Model families (x-axis): Qwen (n=2), Gemma/Gemini (n=2), Llama (n=2), NuExtract (n=2), GPT (n=3), Granite (n=1), GLiNER (n=1)
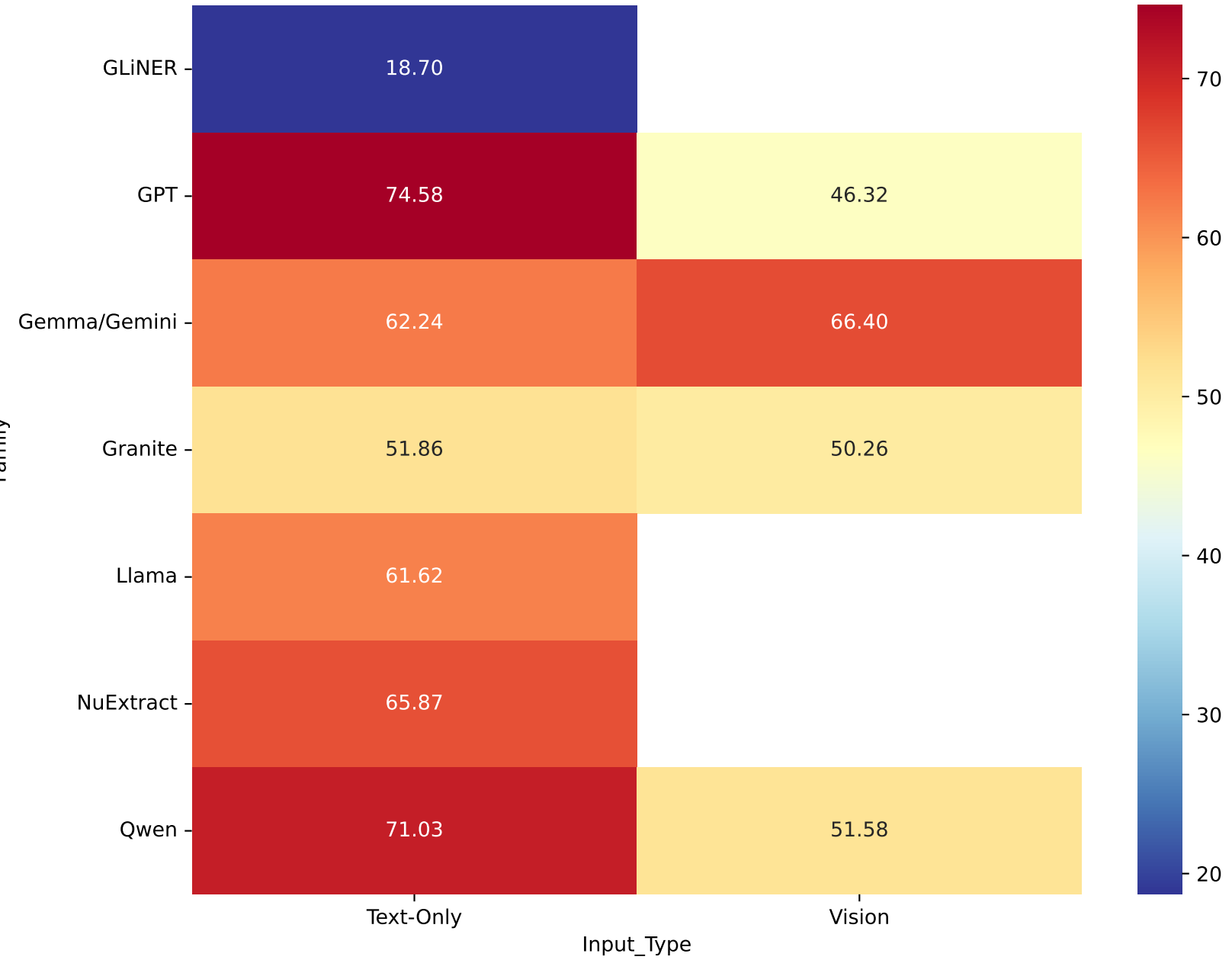
Average Precision

Model Family

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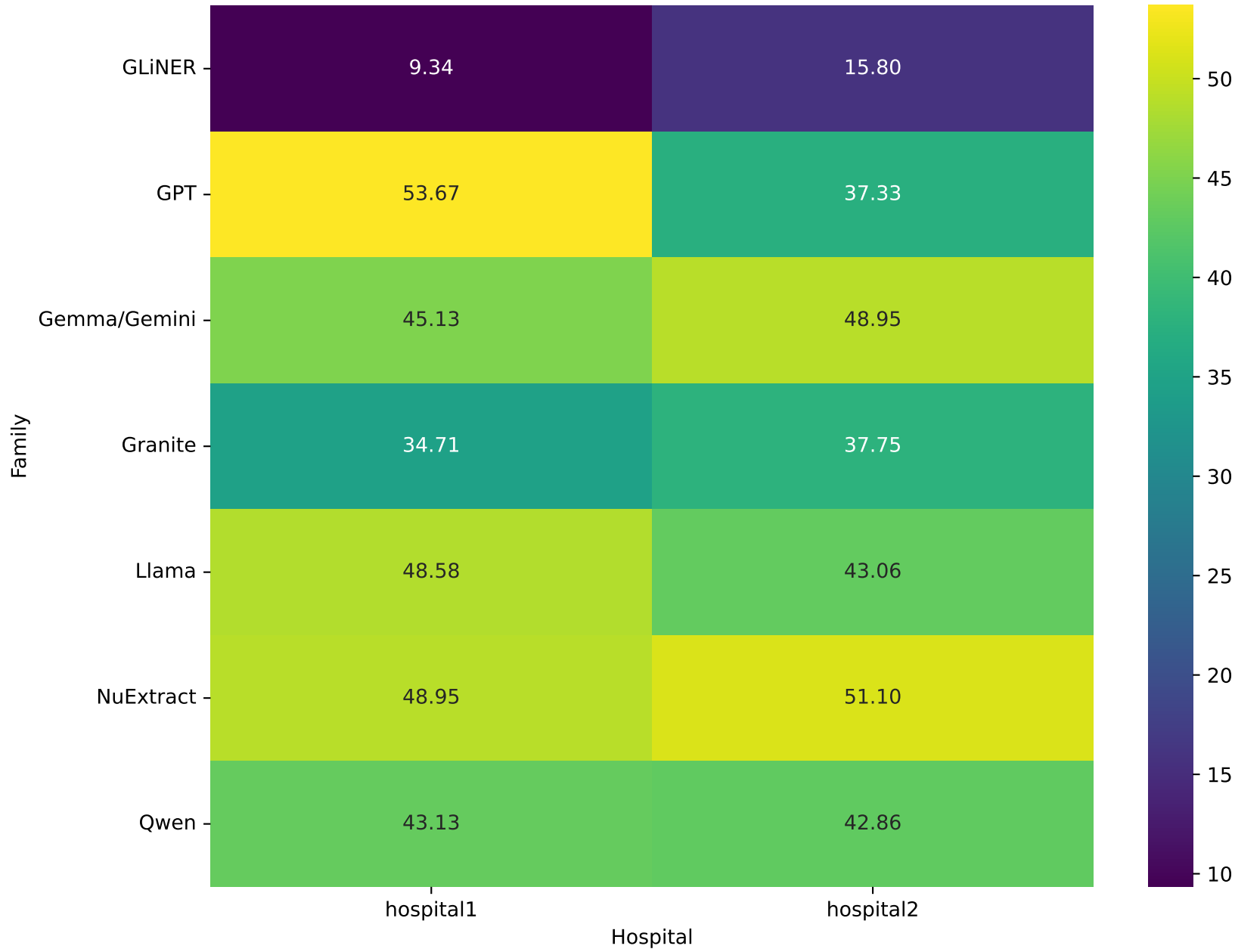**(Averaged across hospitals and input types)**

| Model | Average Recall |
|---|---|
| NuExtract:2B (n=4) | 83.6±5.8 |
| gpt-4o (Image Input, NER Prompt) (n=8) | 71.0±17.7 |
| gpt-4o-mini (Image Input, NER Prompt) (n=8) | 66.2±13.7 |
| llama3.23b (n=4) | 64.8±7.5 |
| gpt-4.1-nano (Image Input, NER Prompt) (n=8) | 64.5±14.8 |
| qwen31.7b (n=4) | 63.7±5.8 |
| gemma34b (Image Input, NER Prompt) (n=8) | 59.4±4.8 |
| NuExtract:0.5B (n=4) | 56.4±10.2 |
| granite3.2 (Image Input, NER Prompt) (n=8) | 53.9±10.1 |
| llama3.21b (n=4) | 53.0±15.0 |
| gemma31b (n=4) | 49.6±10.8 |
| qwen2.5vl3b (Image Input, NER Prompt) (n=8) | 43.6±18.9 |
| GLiNER:NuNerZero (NER Prompt) (n=4) | 12.6±4.3 |

Legend: Qwen, GPT, Gemma/Gemini, Llama, Granite, NuExtract, Given Image Input, Used NER Prompt

**Average Recall by Model Family**

Model Family: NuExtract (n=2), GPT (n=3), Llama (n=2), Gemma/Gemini (n=2), Granite (n=1), Qwen (n=2), GLiNER (n=1)
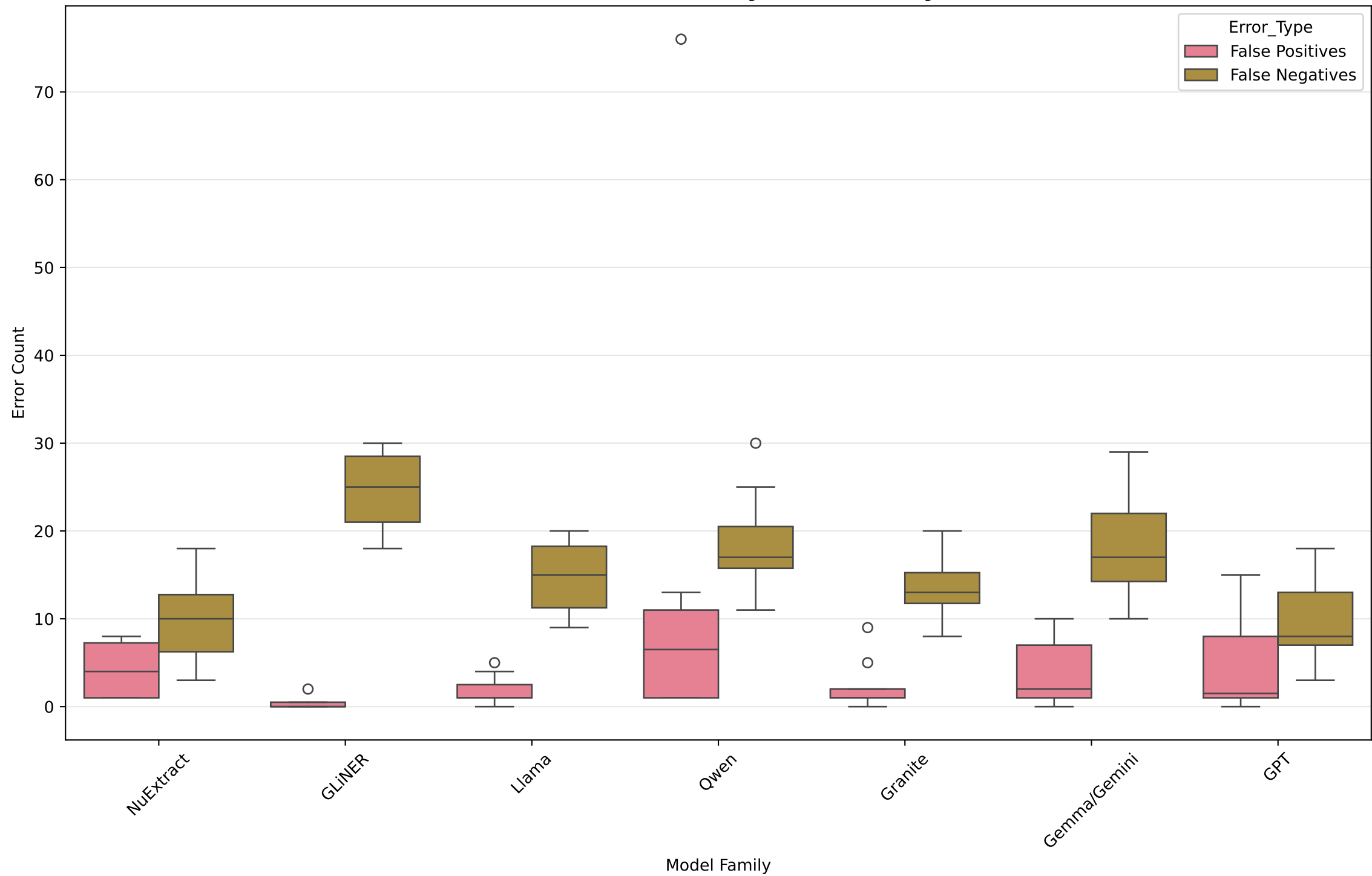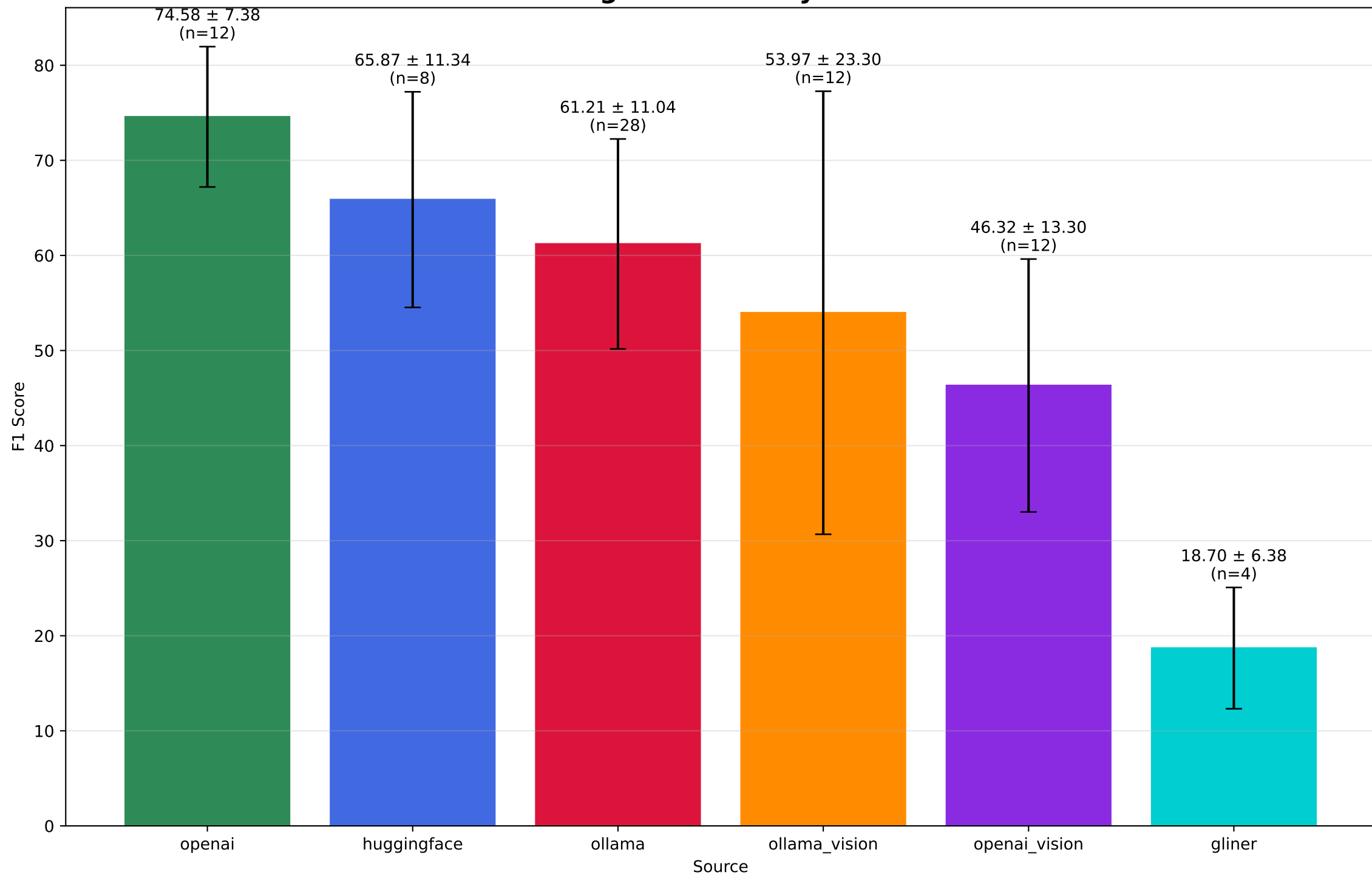
Average F1 Score: Family vs Input Type

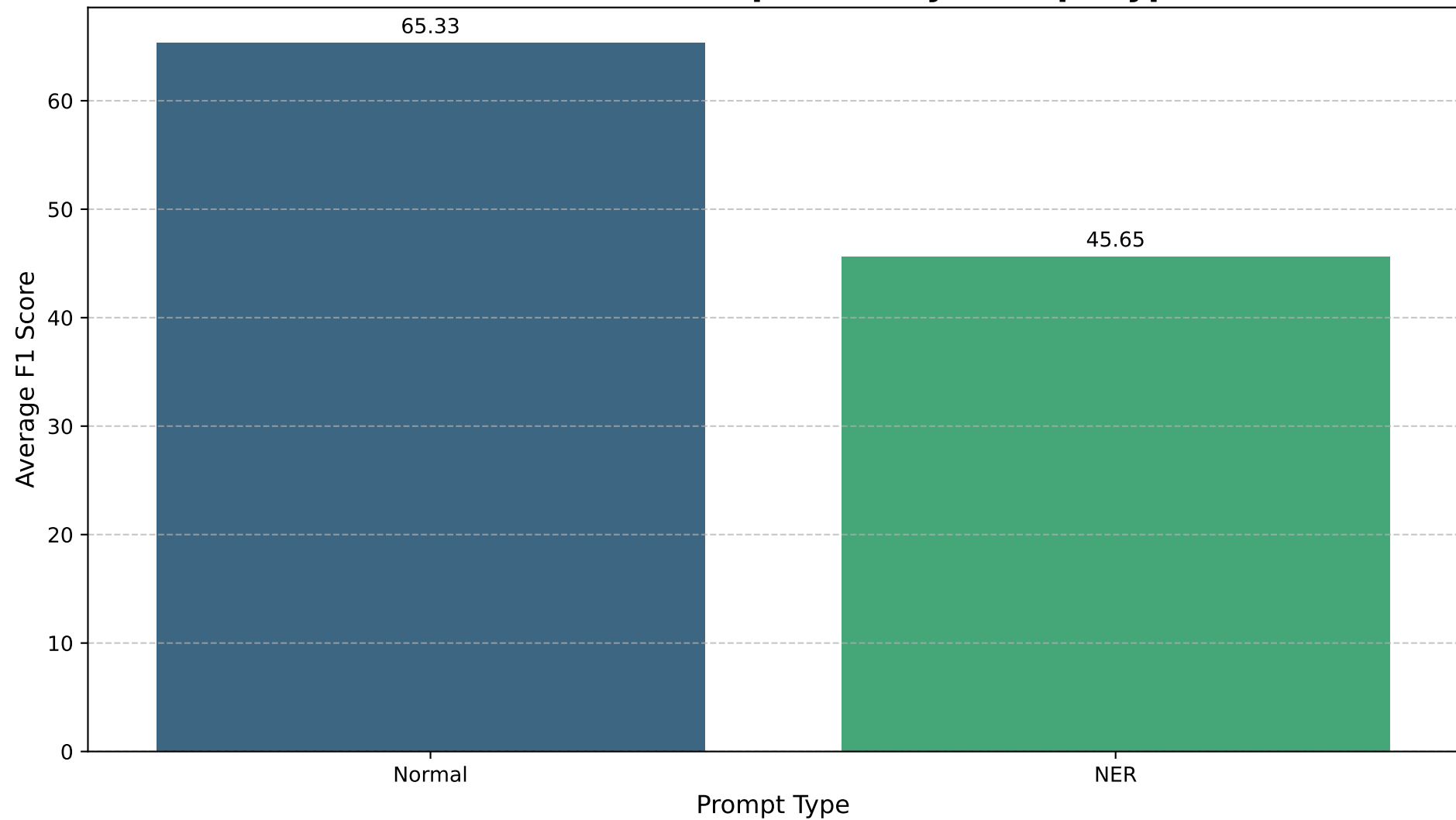Average Accuracy: Family vs Hospital

**Error Distribution by Model Family**

**Average F1 Score by Source**

**Overall F1 Score Comparison by Prompt Type**

Impact of NER Prompt on F1 Score