

Model Pair	Nemenyi p-value	Nemenyi Sig (p<0.05)	Wilcoxon p-value	Wilcoxon Sig (p<0.05)
Claude 3.7 Sonnet vs Llama 3.1 8B	0.0000	✓	0.0000	✓
Claude 3.7 Sonnet vs Llama 3.3 70B	0.0000	✓	0.0000	✓
Claude 3.7 Sonnet vs Minstral 8B	0.0306	✓	0.0000	✓
Claude 3.7 Sonnet vs Qwen 2.5 72B	0.0166	✓	0.0000	✓
Llama 3.1 8B vs Llama 3.3 70B	0.2659		0.0000	✓
Llama 3.1 8B vs Minstral 8B	0.0000	✓	0.0000	✓
Llama 3.1 8B vs Qwen 2.5 72B	0.0000	✓	0.0000	✓
Llama 3.3 70B vs Minstral 8B	0.0166	✓	0.0000	✓
Llama 3.3 70B vs Qwen 2.5 72B	0.0306	✓	0.0000	✓
Minstral 8B vs Qwen 2.5 72B	0.9996		1.0000	