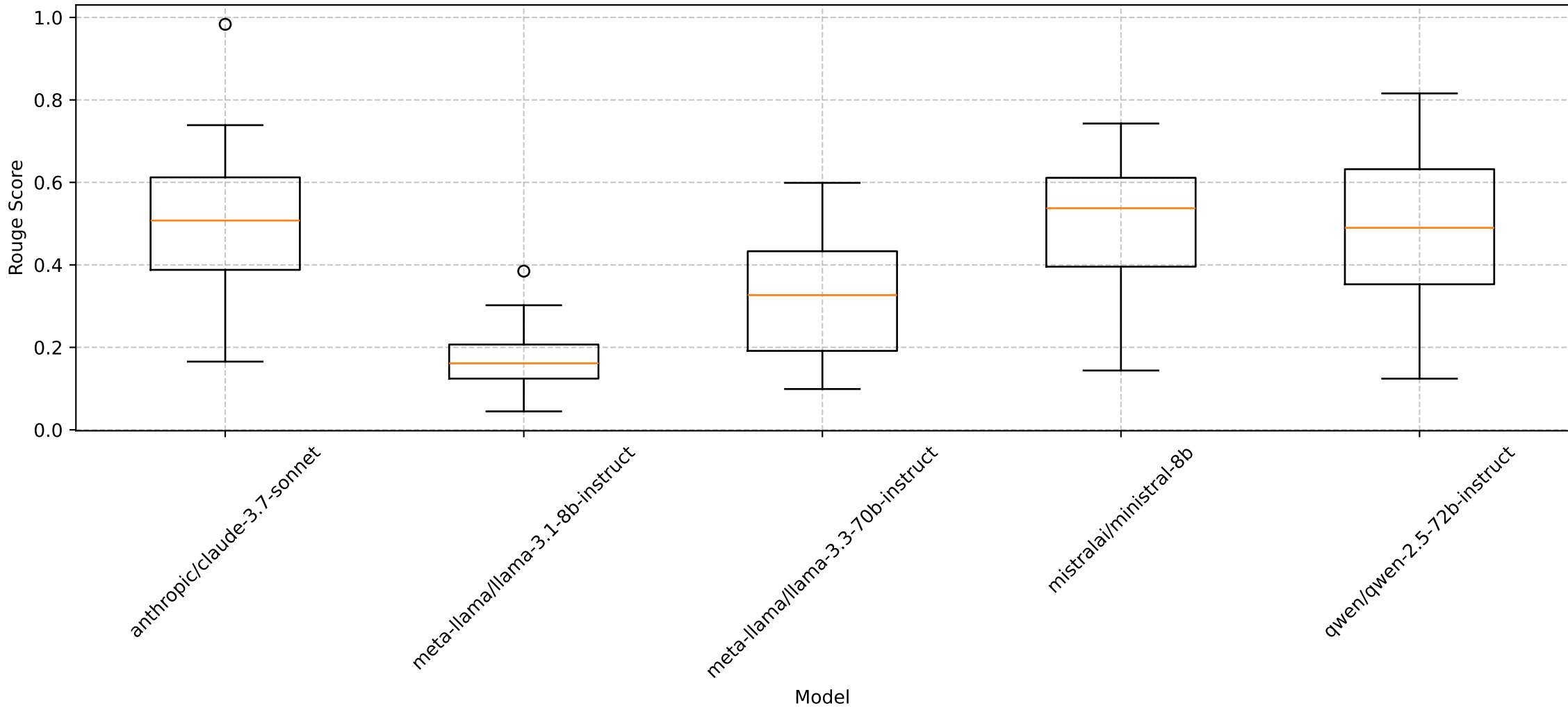
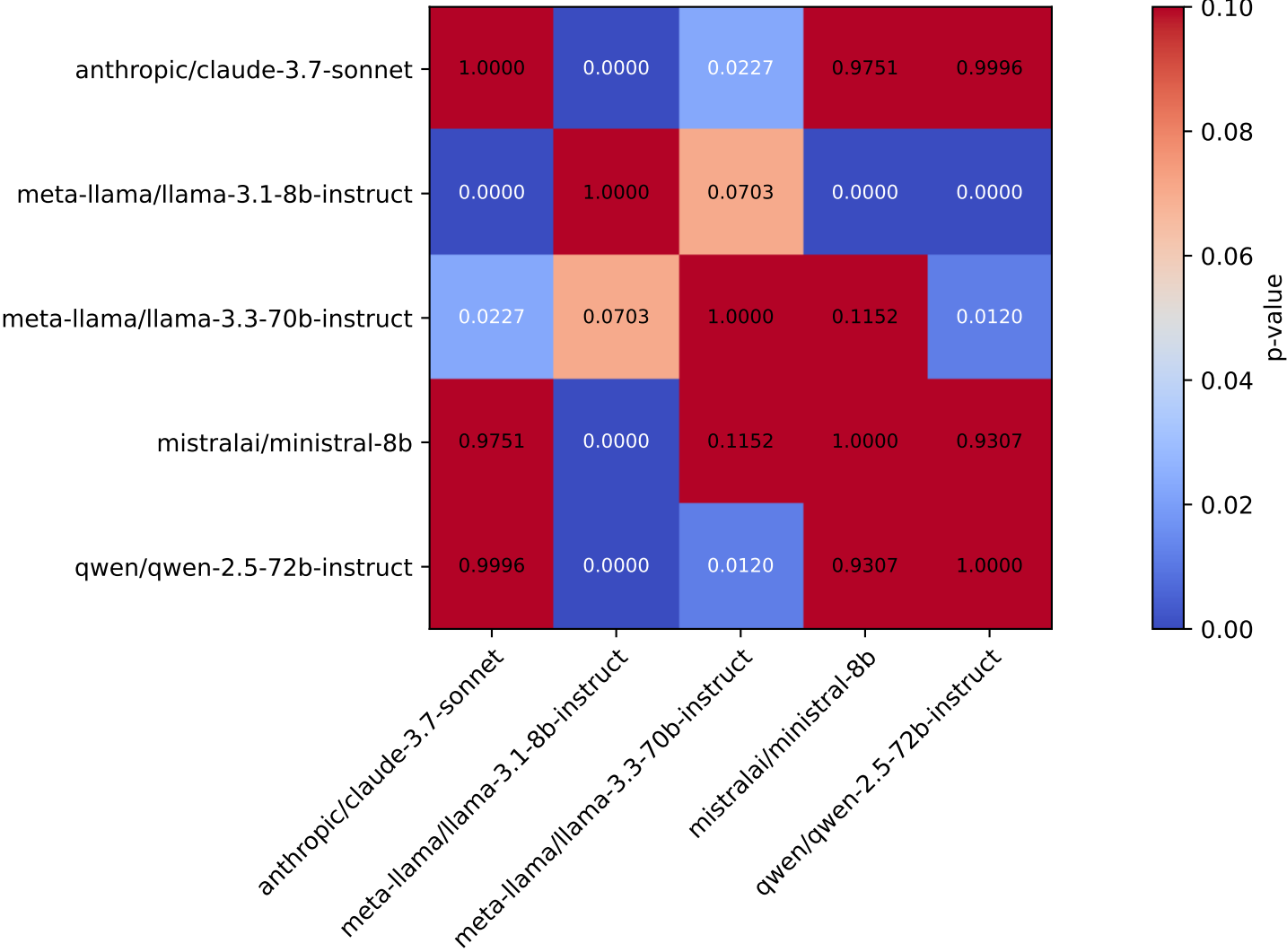


Comparison of rouge_score across models
Friedman $\chi^2 = 49.1200$, $p = 0.0000$ (significant)



Post-hoc Nemenyi Test p-values



Pairwise Wilcoxon Signed-Rank Tests

Comparison	Statistic	p-value	Significant (p<0.05)
anthropic/claude-3.7-sonnet vs meta-llama/llama-3.1-8b-instruct	0.0000	0.0000	✓
anthropic/claude-3.7-sonnet vs meta-llama/llama-3.3-70b-instruct	14.0000	0.0002	✓
anthropic/claude-3.7-sonnet vs mistralai/ministral-8b	92.0000	0.6477	
anthropic/claude-3.7-sonnet vs qwen/qwen-2.5-72b-instruct	92.0000	0.6477	
meta-llama/llama-3.1-8b-instruct vs meta-llama/llama-3.3-70b-instruct	18.0000	0.0005	✓
meta-llama/llama-3.1-8b-instruct vs mistralai/ministral-8b	0.0000	0.0000	✓
meta-llama/llama-3.1-8b-instruct vs qwen/qwen-2.5-72b-instruct	0.0000	0.0000	✓
meta-llama/llama-3.3-70b-instruct vs mistralai/ministral-8b	11.0000	0.0001	✓
meta-llama/llama-3.3-70b-instruct vs qwen/qwen-2.5-72b-instruct	20.0000	0.0007	✓
mistralai/ministral-8b vs qwen/qwen-2.5-72b-instruct	85.0000	0.4749	