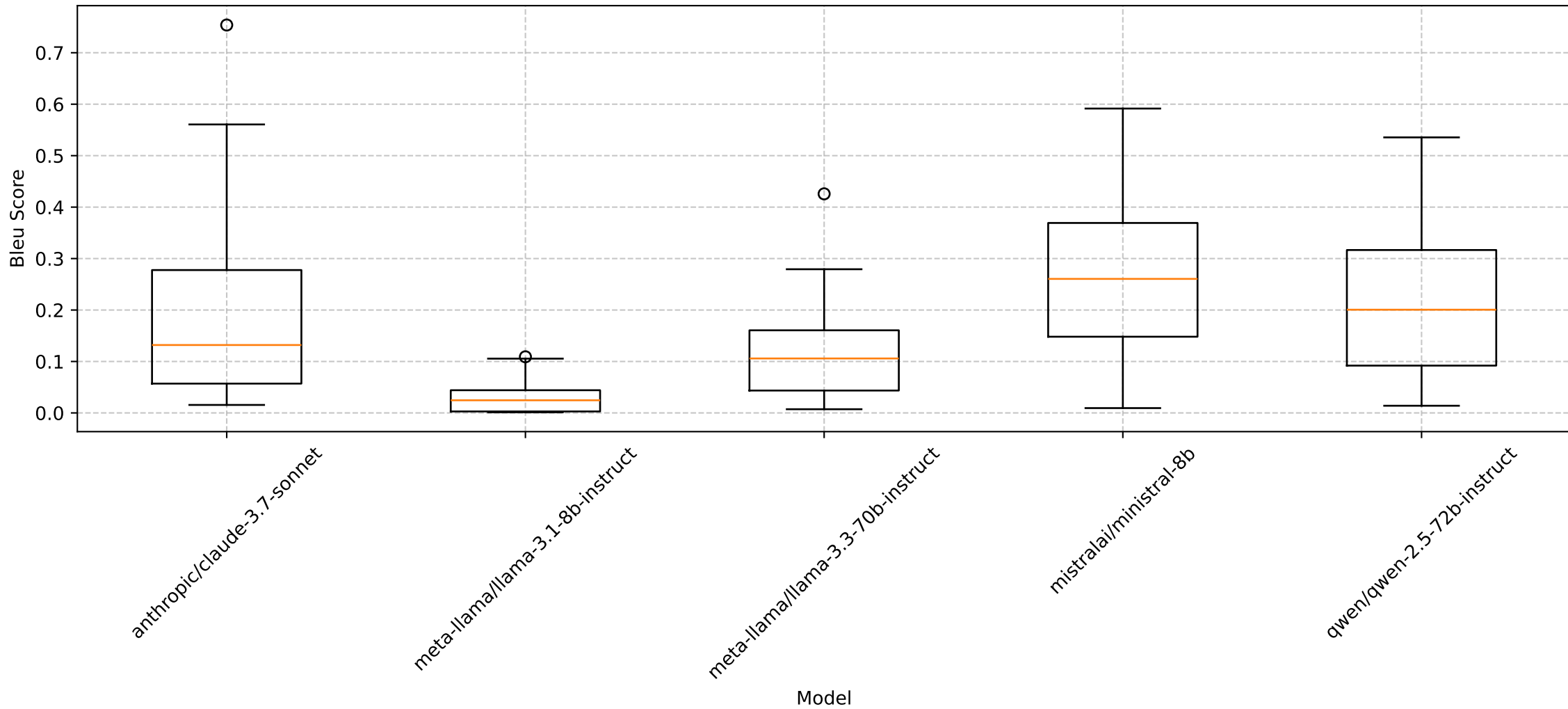
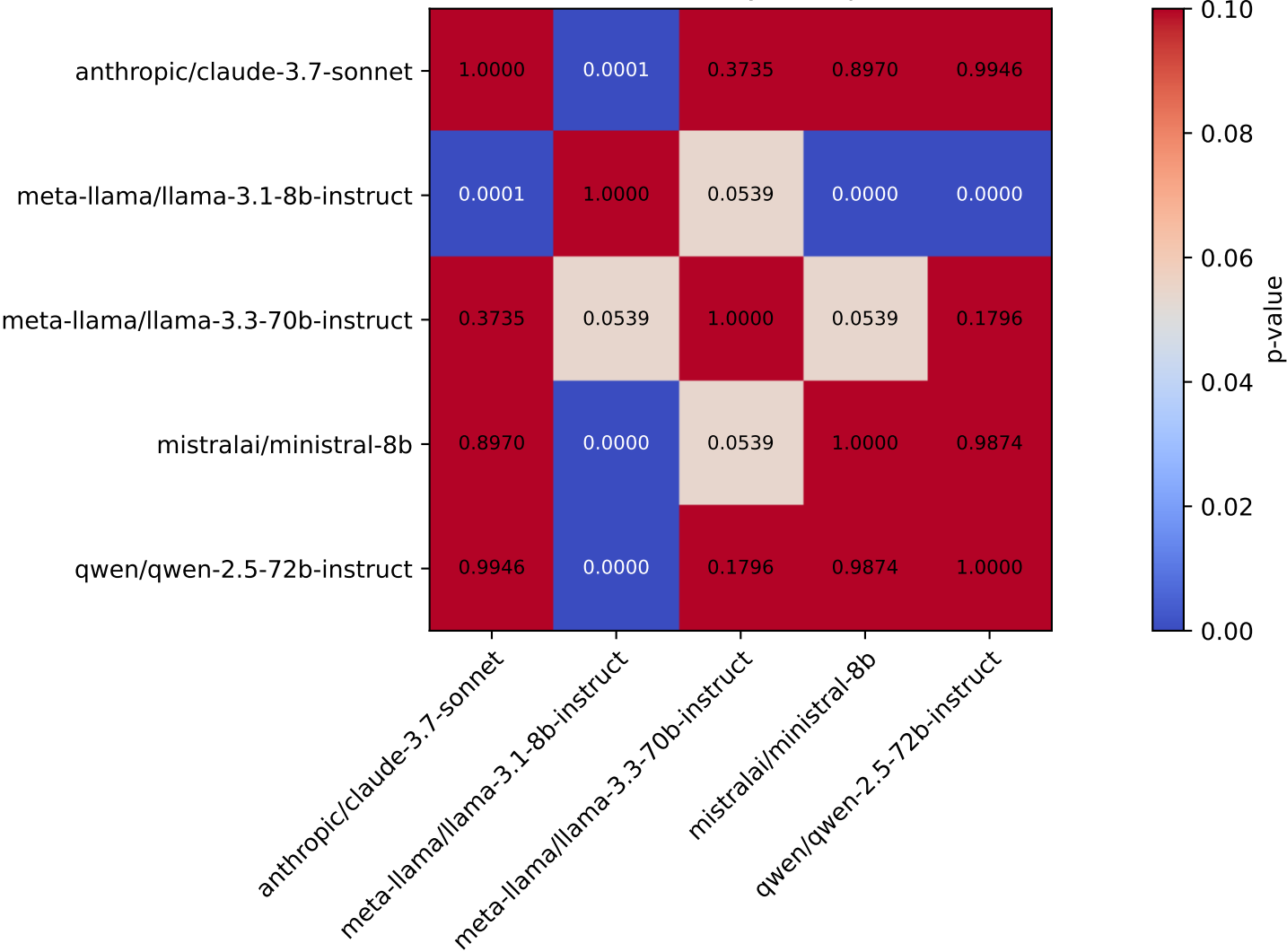


Comparison of bleu\_score across models  
Friedman  $\chi^2 = 38.9200$ ,  $p = 0.0000$  (significant)



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	5.0000	0.0000	✓
anthropic/claude-3.7-sonnet vs meta-llama/llama-3.3-70b-instruct	67.0000	0.1650	
anthropic/claude-3.7-sonnet vs mistralai/ministral-8b	64.0000	0.1327	
anthropic/claude-3.7-sonnet vs qwen/qwen-2.5-72b-instruct	102.0000	0.9273	
meta-llama/llama-3.1-8b-instruct vs meta-llama/llama-3.3-70b-instruct	11.0000	0.0001	✓
meta-llama/llama-3.1-8b-instruct vs mistralai/ministral-8b	2.0000	0.0000	✓
meta-llama/llama-3.1-8b-instruct vs qwen/qwen-2.5-72b-instruct	2.0000	0.0000	✓
meta-llama/llama-3.3-70b-instruct vs mistralai/ministral-8b	19.0000	0.0006	✓
meta-llama/llama-3.3-70b-instruct vs qwen/qwen-2.5-72b-instruct	31.0000	0.0042	✓
mistralai/ministral-8b vs qwen/qwen-2.5-72b-instruct	90.0000	0.5958	