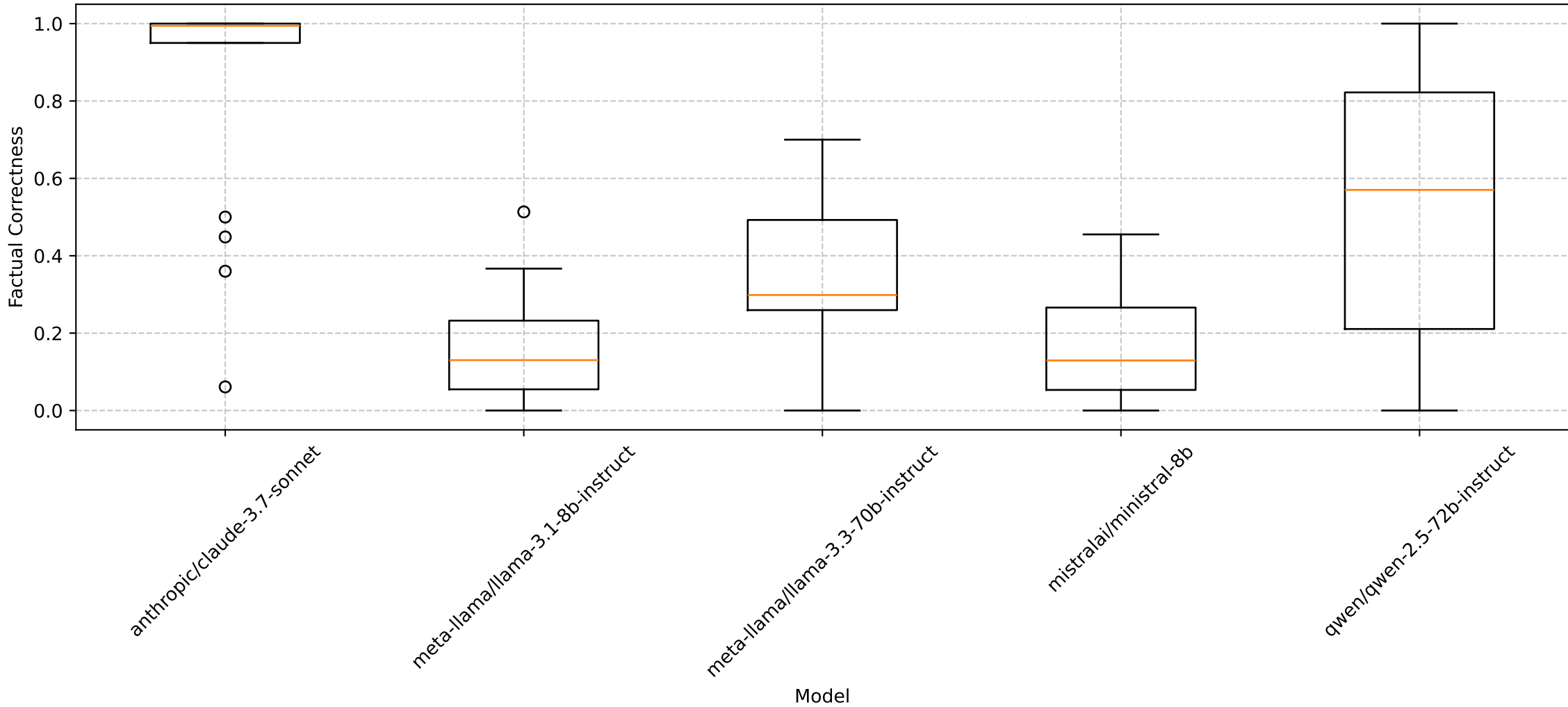
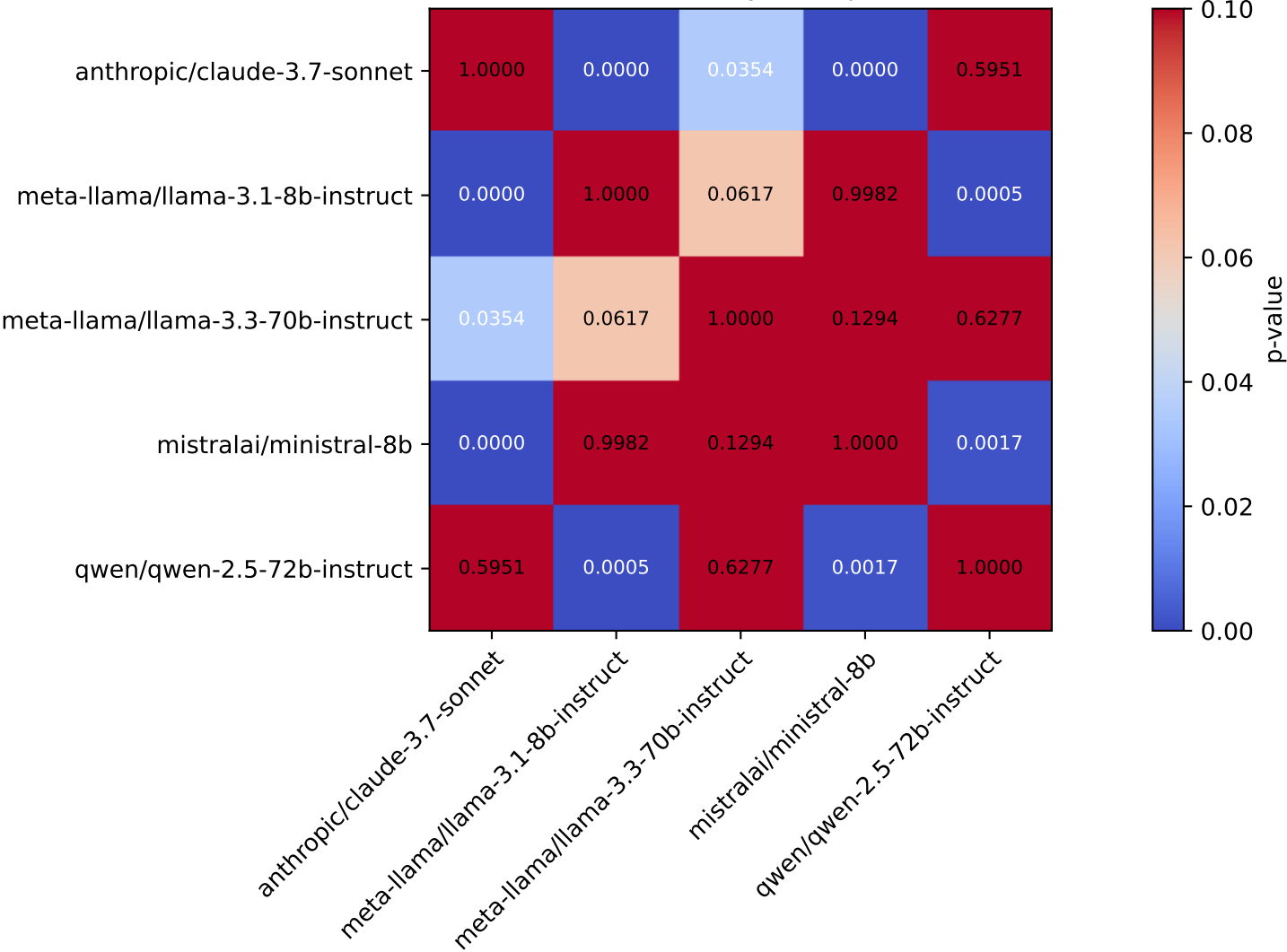


Comparison of factual correctness across models
Friedman $\chi^2 = 45.1429$, $p = 0.0000$ (significant)



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	2.0000	0.0000	✓
anthropic/claude-3.7-sonnet vs meta-llama/llama-3.3-70b-instruct	6.0000	0.0000	✓
anthropic/claude-3.7-sonnet vs mistralai/ministral-8b	2.0000	0.0000	✓
anthropic/claude-3.7-sonnet vs qwen/qwen-2.5-72b-instruct	25.0000	0.0017	✓
meta-llama/llama-3.1-8b-instruct vs meta-llama/llama-3.3-70b-instruct	22.0000	0.0010	✓
meta-llama/llama-3.1-8b-instruct vs mistralai/ministral-8b	98.0000	0.8124	
meta-llama/llama-3.1-8b-instruct vs qwen/qwen-2.5-72b-instruct	13.0000	0.0002	✓
meta-llama/llama-3.3-70b-instruct vs mistralai/ministral-8b	22.0000	0.0010	✓
meta-llama/llama-3.3-70b-instruct vs qwen/qwen-2.5-72b-instruct	37.0000	0.0196	✓
mistralai/ministral-8b vs qwen/qwen-2.5-72b-instruct	13.0000	0.0002	✓