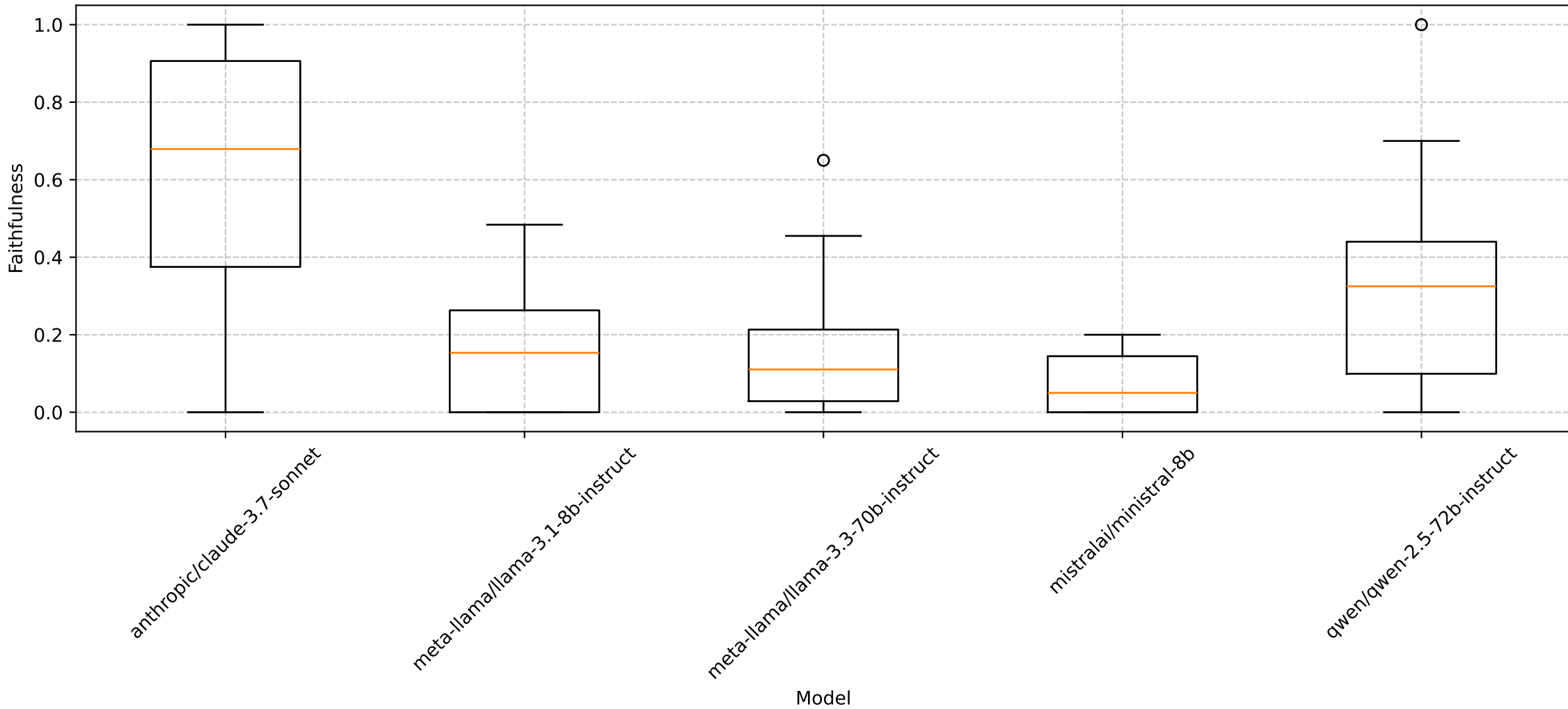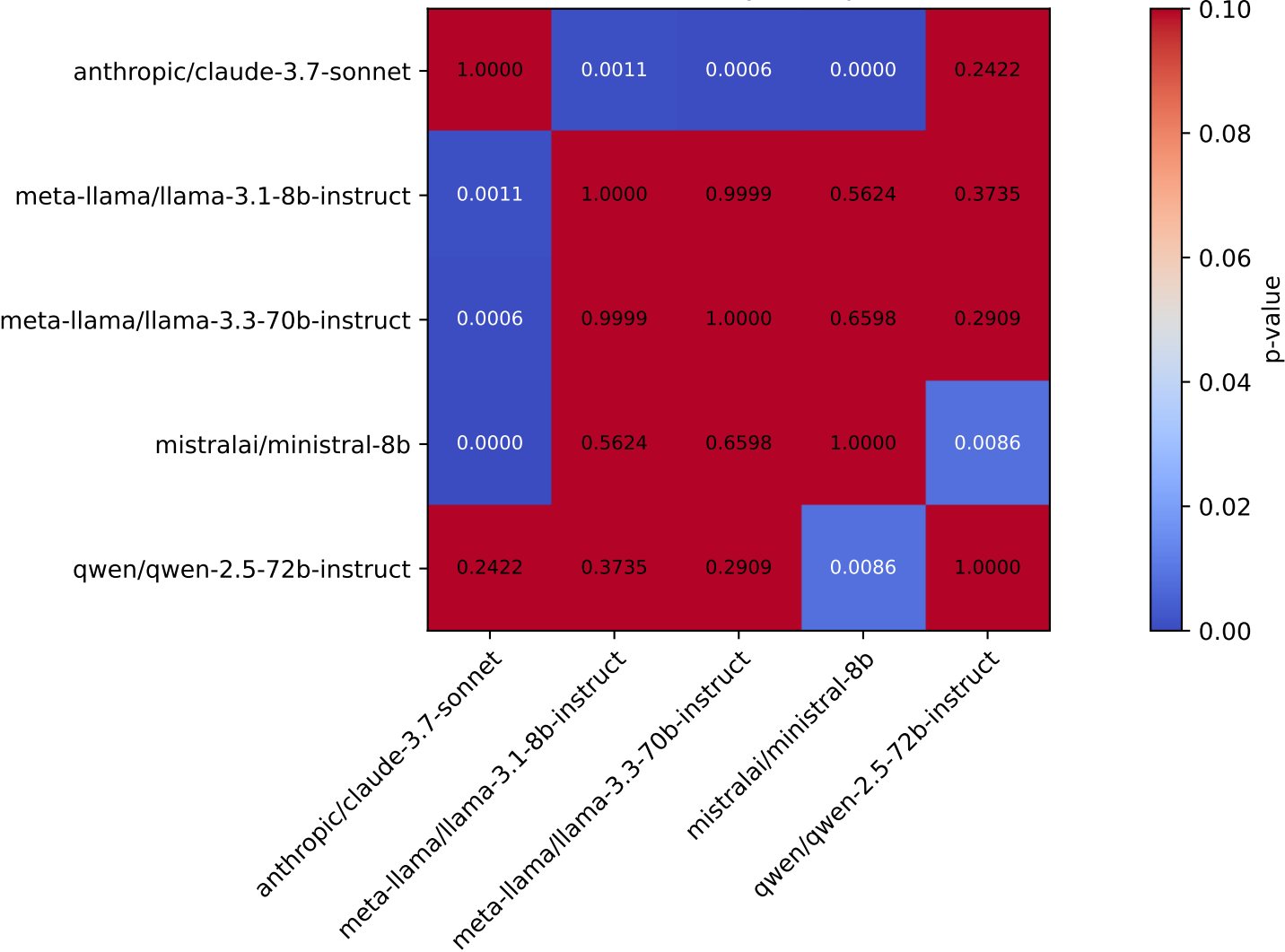Comparison of faithfulness across models
Friedman $\chi^2$ = 37.8674, p = 0.0000 (significant)

## Post-hoc Nemenyi Test p-values



|  | anthropic/claude-3.7-sonnet | meta-llama/llama-3.1-8b-instruct | meta-llama/llama-3.3-70b-instruct | mistralai/ministral-8b | qwen/qwen-2.5-72b-instruct |
|---|---|---|---|---|---|
| anthropic/claude-3.7-sonnet | 1.0000 | 0.0011 | 0.0006 | 0.0000 | 0.2422 |
| meta-llama/llama-3.1-8b-instruct | 0.0011 | 1.0000 | 0.9999 | 0.5624 | 0.3735 |
| meta-llama/llama-3.3-70b-instruct | 0.0006 | 0.9999 | 1.0000 | 0.6598 | 0.2909 |
| mistralai/ministral-8b | 0.0000 | 0.5624 | 0.6598 | 1.0000 | 0.0086 |
| qwen/qwen-2.5-72b-instruct | 0.2422 | 0.3735 | 0.2909 | 0.0086 | 1.0000 |

## 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 | 4.0000 | 0.0002 | ✓ |
| anthropic/claude-3.7-sonnet vs meta-llama/llama-3.3-70 | 6.0000 | 0.0005 | ✓ |
| anthropic/claude-3.7-sonnet vs mistralai/ministral | 0.0000 | 0.0002 | ✓ |
| anthropic/claude-3.7-sonnet vs qwen/qwen-2.5-72b-i | 15.0000 | 0.0021 | ✓ |
| llama/llama-3.1-8b-instruct vs meta-llama/llama-3.3- | 73.0000 | 0.5861 |  |
| meta-llama/llama-3.1-8b-instruct vs mistralai/minist | 17.0000 | 0.0083 | ✓ |
| llama/llama-3.1-8b-instruct vs qwen/qwen-2.5-72b | 25.0000 | 0.0048 | ✓ |
| ta-llama/llama-3.3-70b-instruct vs mistralai/minist | 28.0000 | 0.0385 | ✓ |
| llama/llama-3.3-70b-instruct vs qwen/qwen-2.5-72 | 11.0000 | 0.0019 | ✓ |
| mistralai/ministral-8b vs qwen/qwen-2.5-72b-instr | 8.0000 | 0.0007 | ✓ |