# Calibration and Metacognition Metrics (lower is better)



**Expected Calibration Error**

| Model | ECE |
|-------|-----|
| Claude-Sonnet-4 | 0.560 |
| Kimi-K2 | 0.559 |
| DeepSeek-V3.1 | 0.554 |
| Gemini-2.5 | 0.503 |
| GPT-4.1 | 0.472 |
| Qwen3-235B | 0.451 |

**Brier Score**

| Model | Brier Score |
|-------|-------------|
| Claude-Sonnet-4 | 0.558 |
| Kimi-K2 | 0.558 |
| DeepSeek-V3.1 | 0.552 |
| Gemini-2.5 | 0.503 |
| GPT-4.1 | 0.472 |
| Qwen3-235B | 0.452 |

**Wrong@High-Confidence (0.9)**

| Model | Error Rate (%) |
|-------|----------------|
| Claude-Sonnet-4 | 0.0% |
| Kimi-K2 | 0.0% |
| DeepSeek-V3.1 | 0.0% |
| Gemini-2.5 | 0.0% |
| GPT-4.1 | 0.0% |
| Qwen3-235B | 0.0% |