

Top Model Performance Comparison

Memory Efficiency

Loading Speed

Size Efficiency

Inference Speed

Response Quality

- Florence-2-base
- Moondream2-2B
- Qwen2.5-VL-3B
- Qwen2.5-VL-7B

