**LLaMA 70B Inference Performance (M2 Max)**

Throughput (tokens/second) vs Quantization Scheme

| Quantization Scheme | Throughput (tokens/second) |
| --- | --- |
| FP16 | 2 |
| INT8 | 4 |
| INT4 | 8 |
| 3.5-bit | 9 |