

Scalability across Model Sizes (Decode)

Decode Throughput (Tokens/s)

