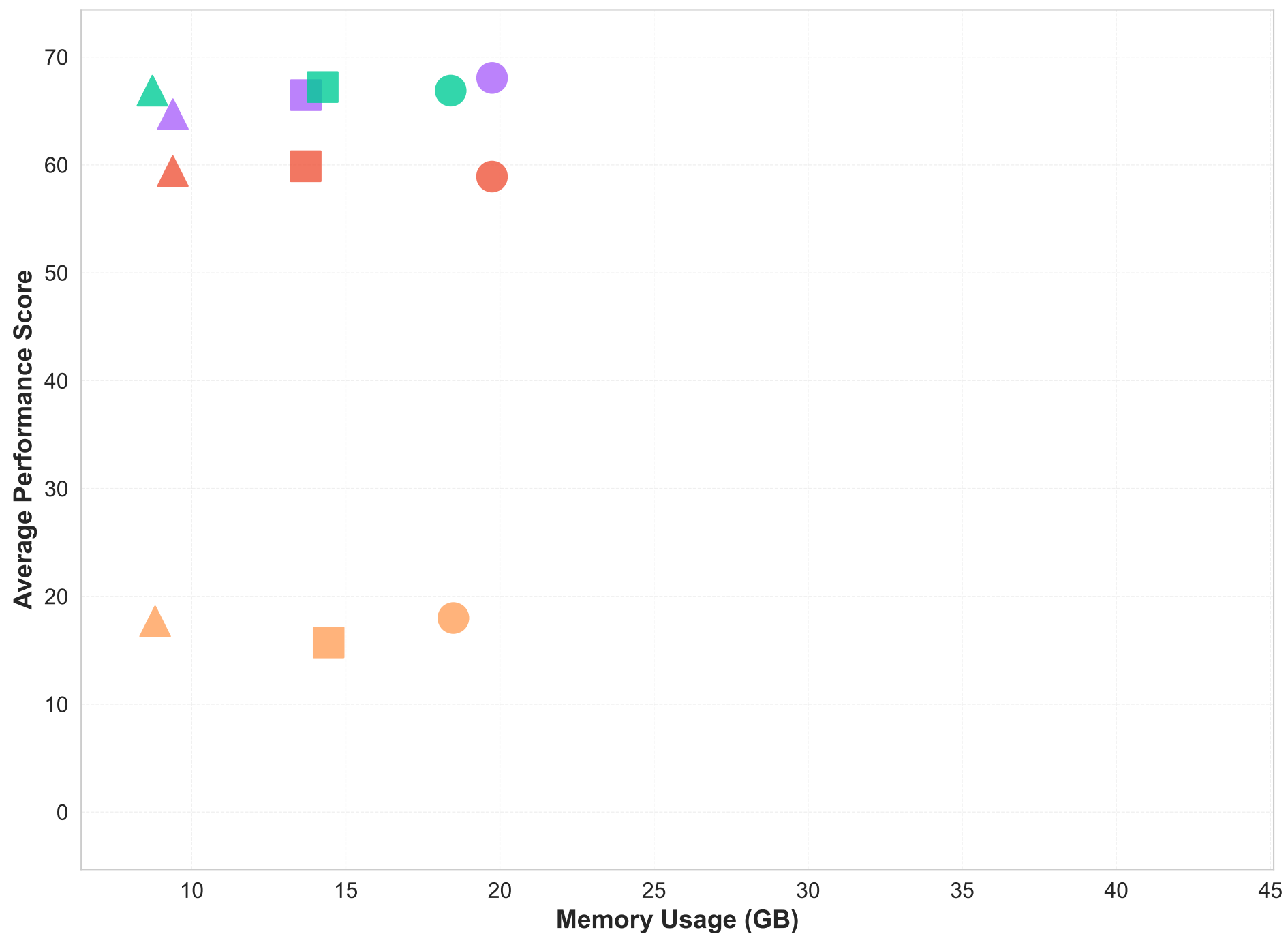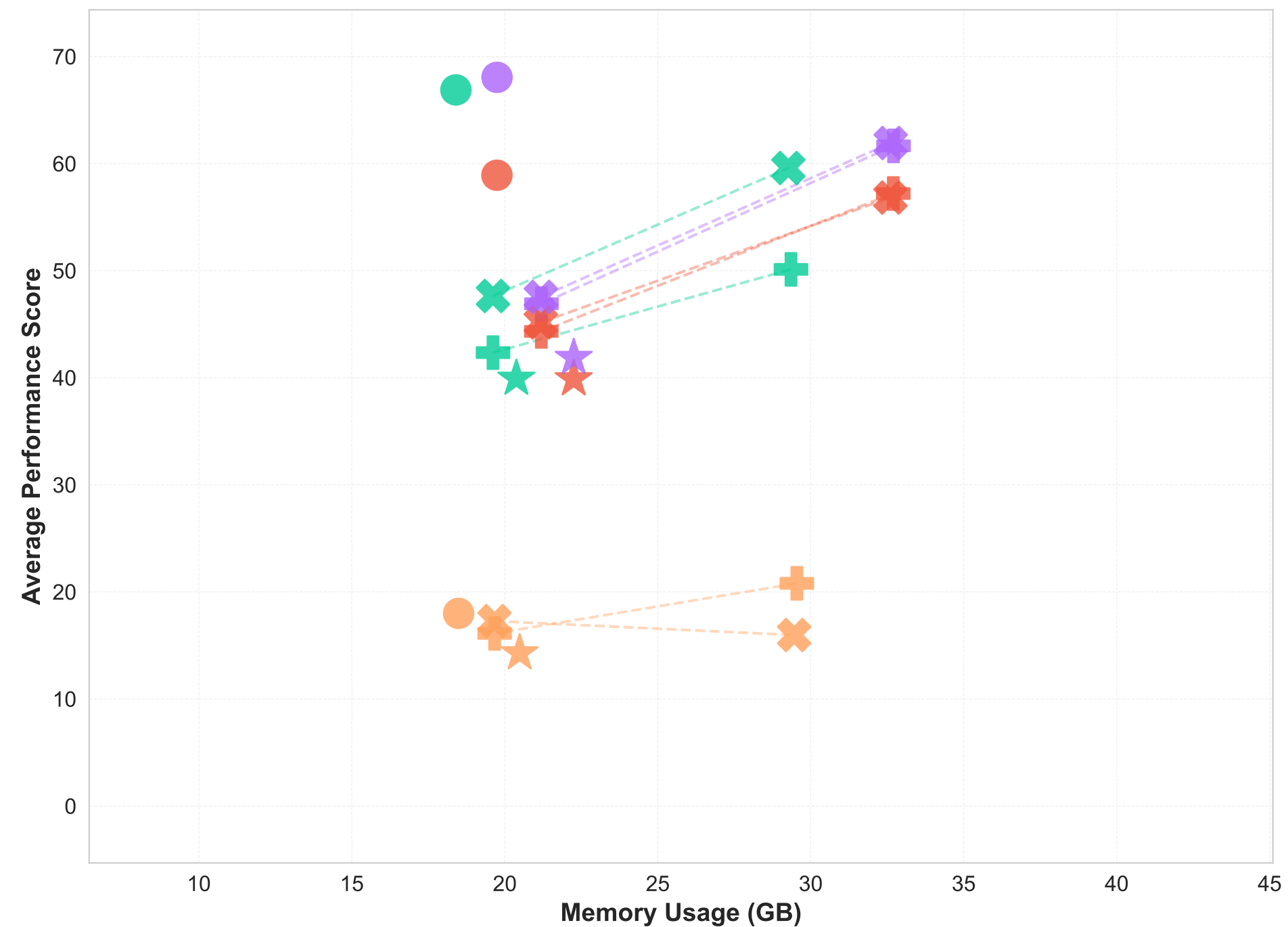Model Comparison: Average Performance vs Memory Usage
HELMET (16K) and LongProc (2K) Benchmarks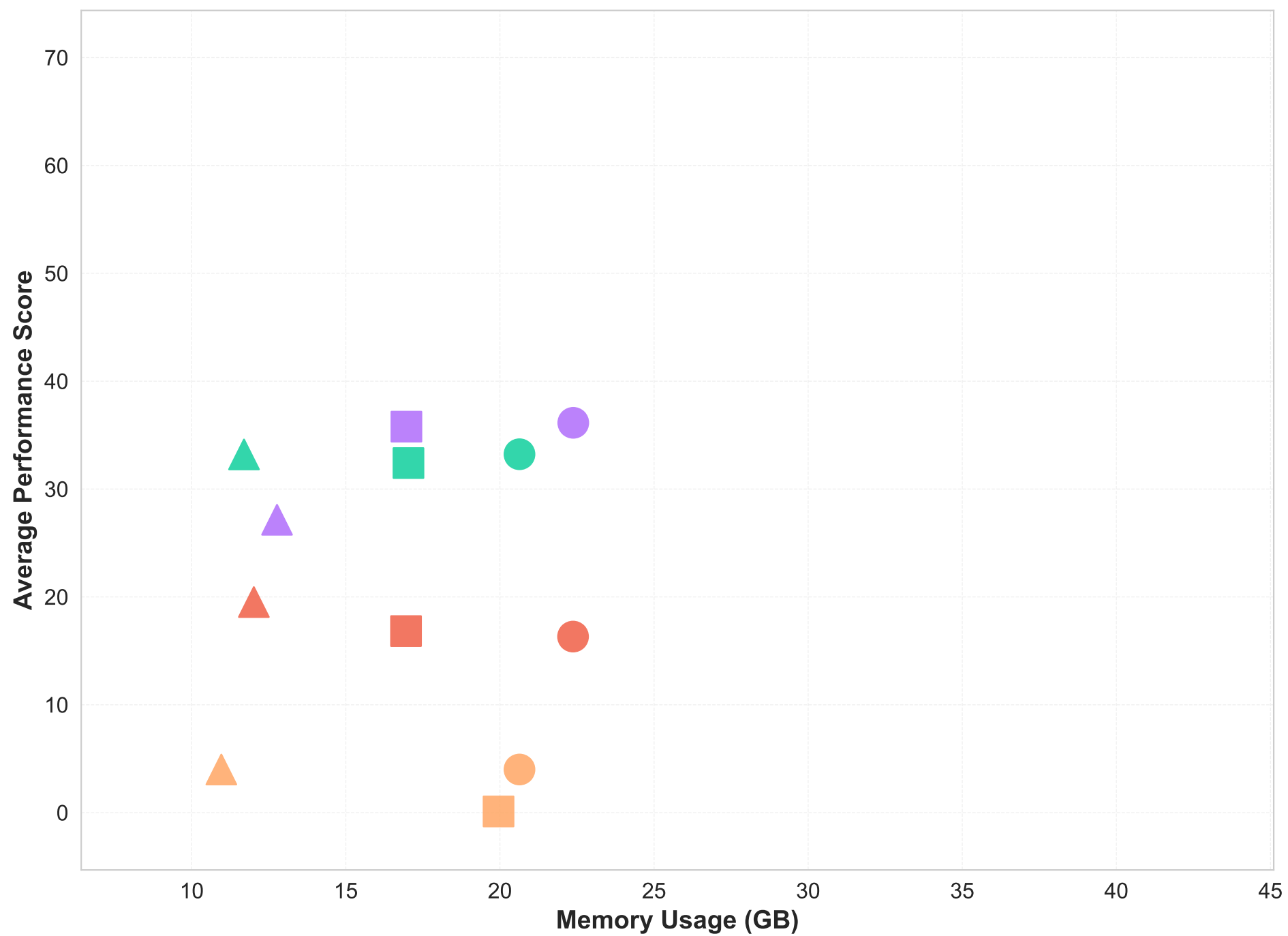