Token Eviction Methods Struggle on the Pareto Frontier of Performance vs Efficiency
HELMET Rerank Task: Memory vs NDCG@10 (16K Context)