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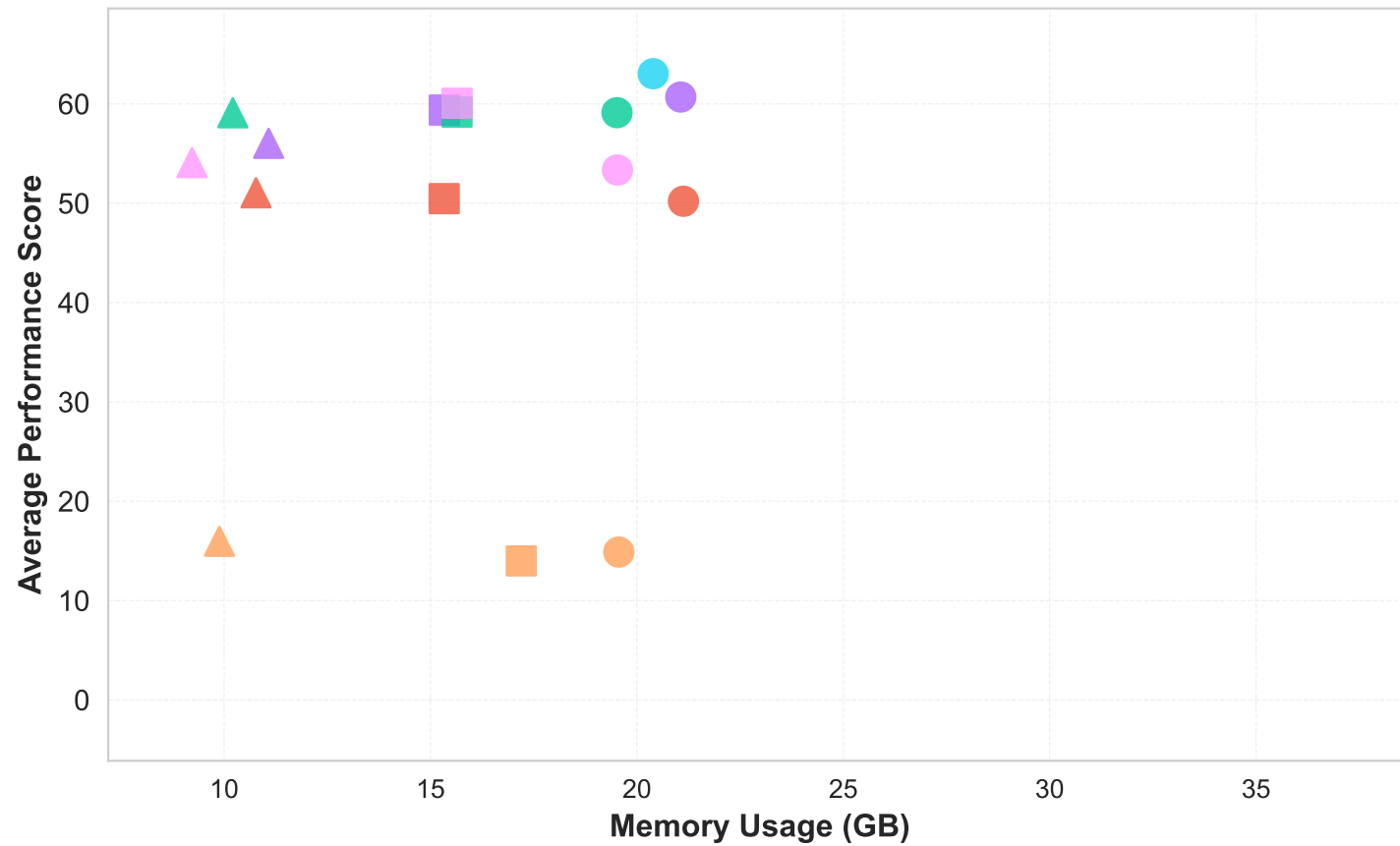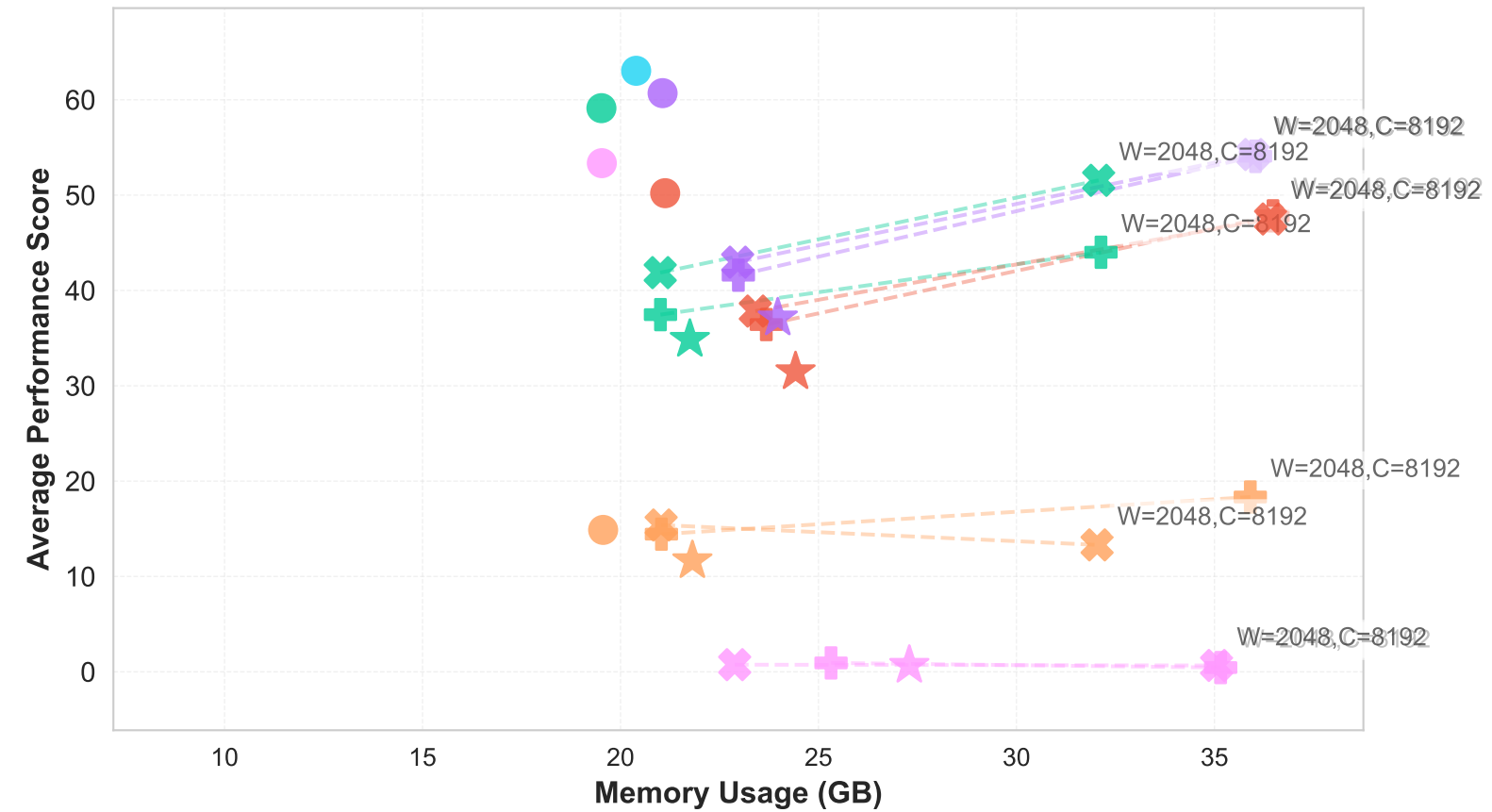**Combined HELMET (16K) and LongProc (2K) Benchmarks**

Quantization Methods

KV Cache Methods

Average Performance Score

Memory Usage (GB)

Models:
- Llama-3.1-8B-Instruct
- Qwen2.5-7B-Instruct
- DeepSeek-R1-Distill-Llama-8B
- DeepSeek-R1-Distill-Qwen-7B
- Qwen3-8B
- Yarn-Qwen3-8B

Techniques:
- INT4
- INT8
- baseline
- pyramidkv
- snapkv
- streamingllm

W=2048,C=8192