

Model

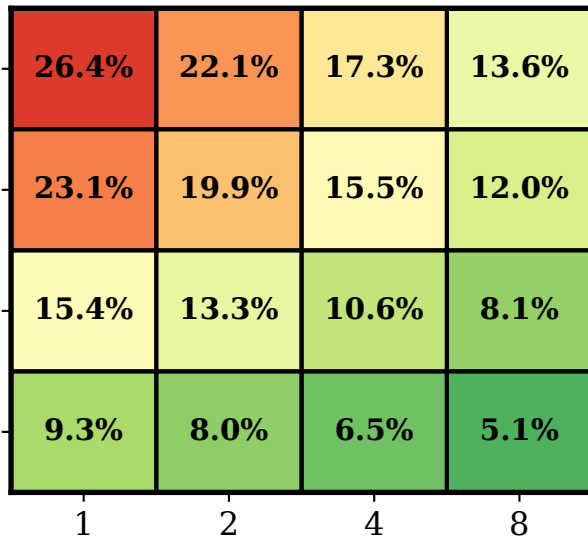
Baseline

Conservative

Moderate

Aggressive

Switch-Base



OLMoE-1B-7B

Qwen1.5-MoE-A2.7B

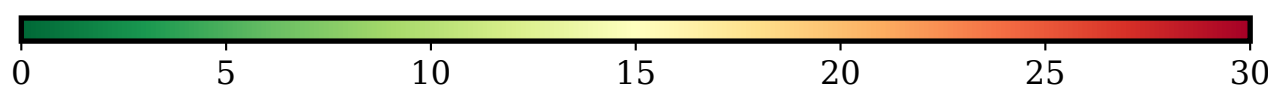
Switch-Large

Batch Size

Batch Size

Batch Size

Batch Size



Relative Overhead (%)