

## Inference peak memory: Baseline vs INT4

