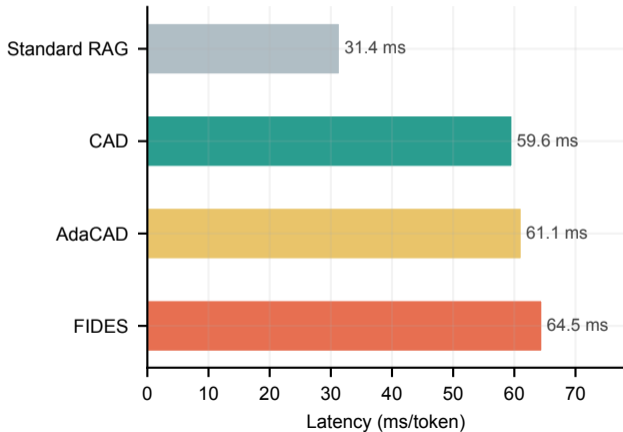


# Per-Token Inference Latency

## LLaMA3-8B



## Qwen3-8B

