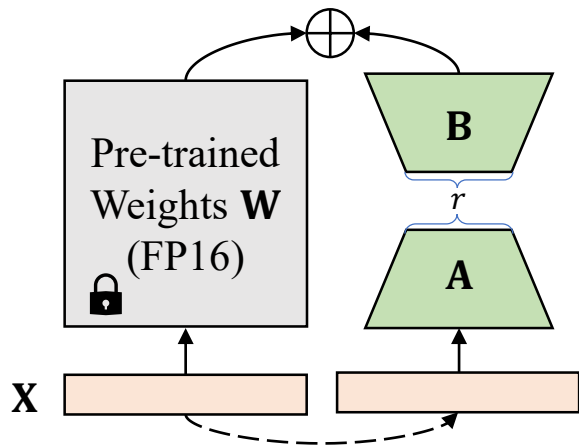
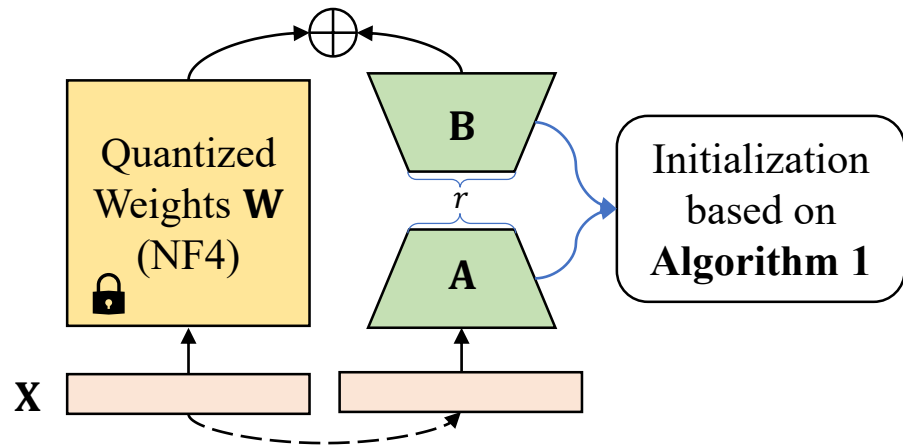


Fine-tuning

Full-precision LoRA



LoRA Fine-tuning-aware Quantization



Inference

