Training Phase

1. Retraining

Training Data



Pretrained LLM

2. Fine-Tuning

Instruction Data



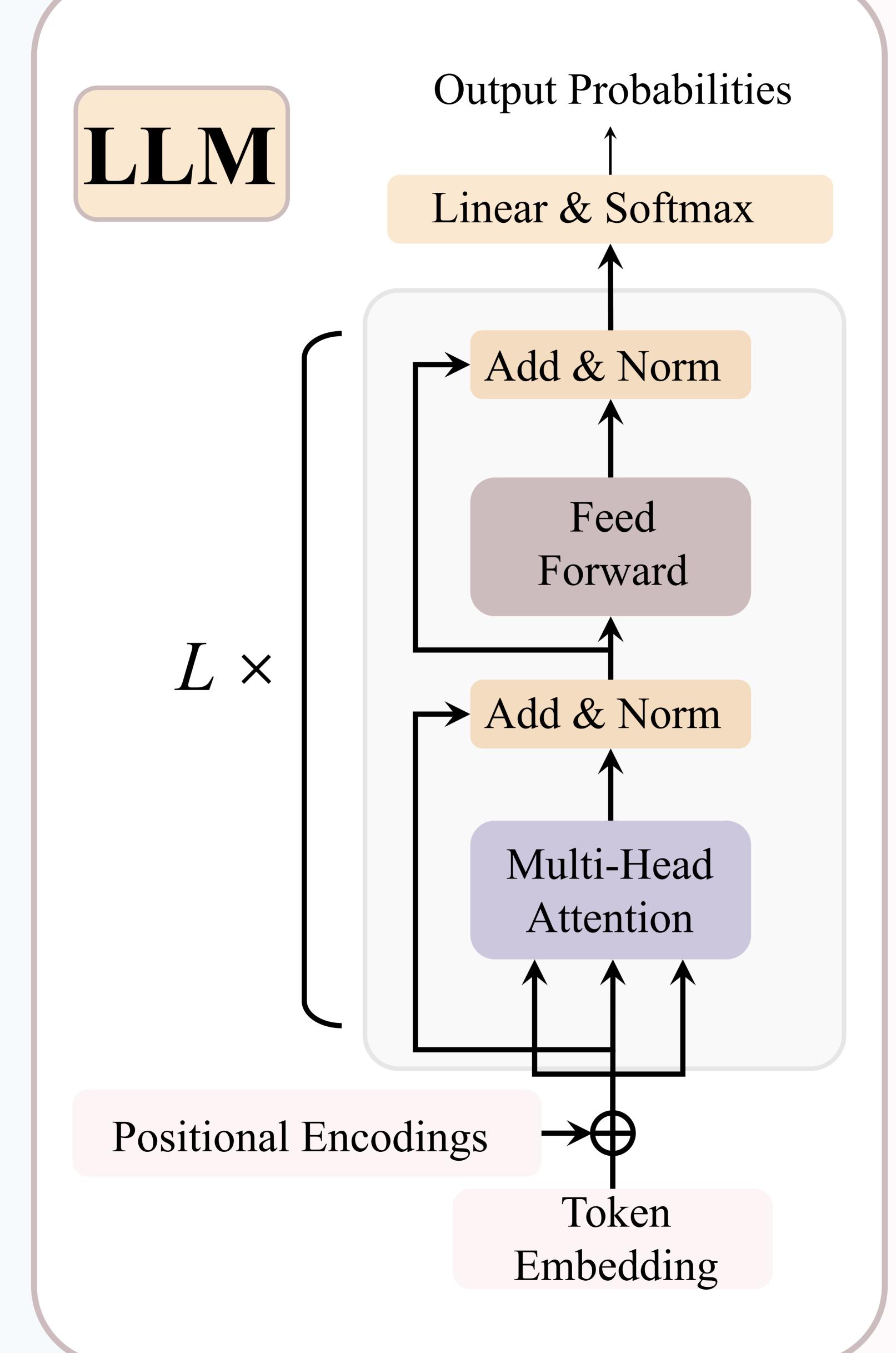
Chat LLM



RL Feedback

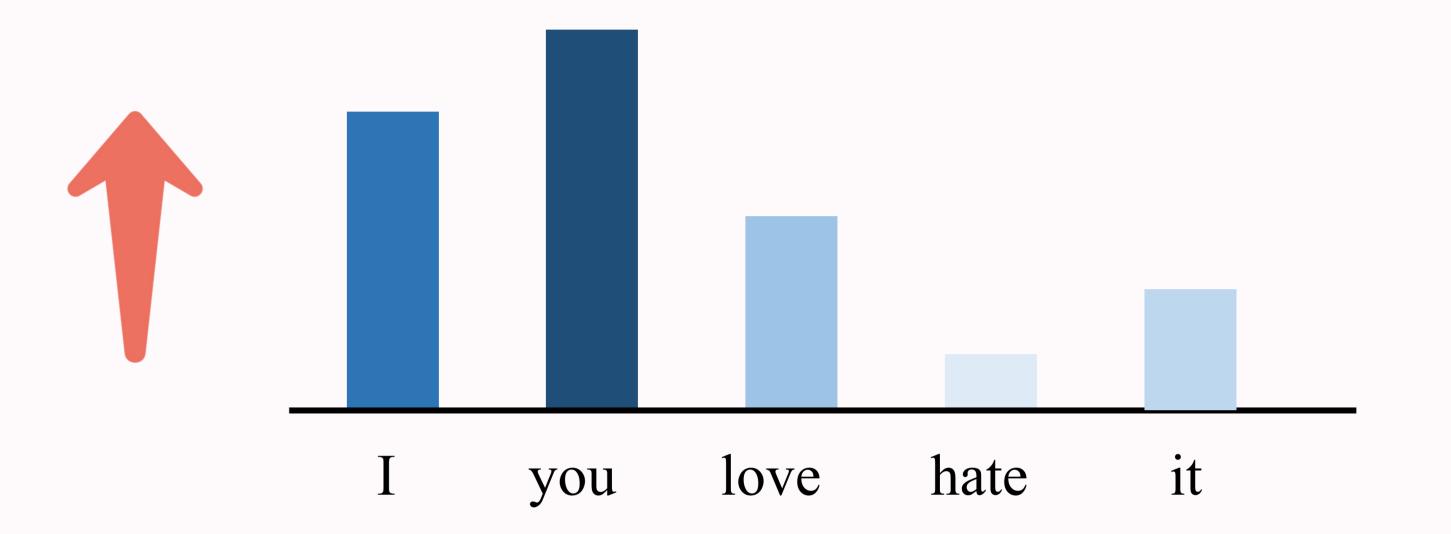


Aligned LLIV

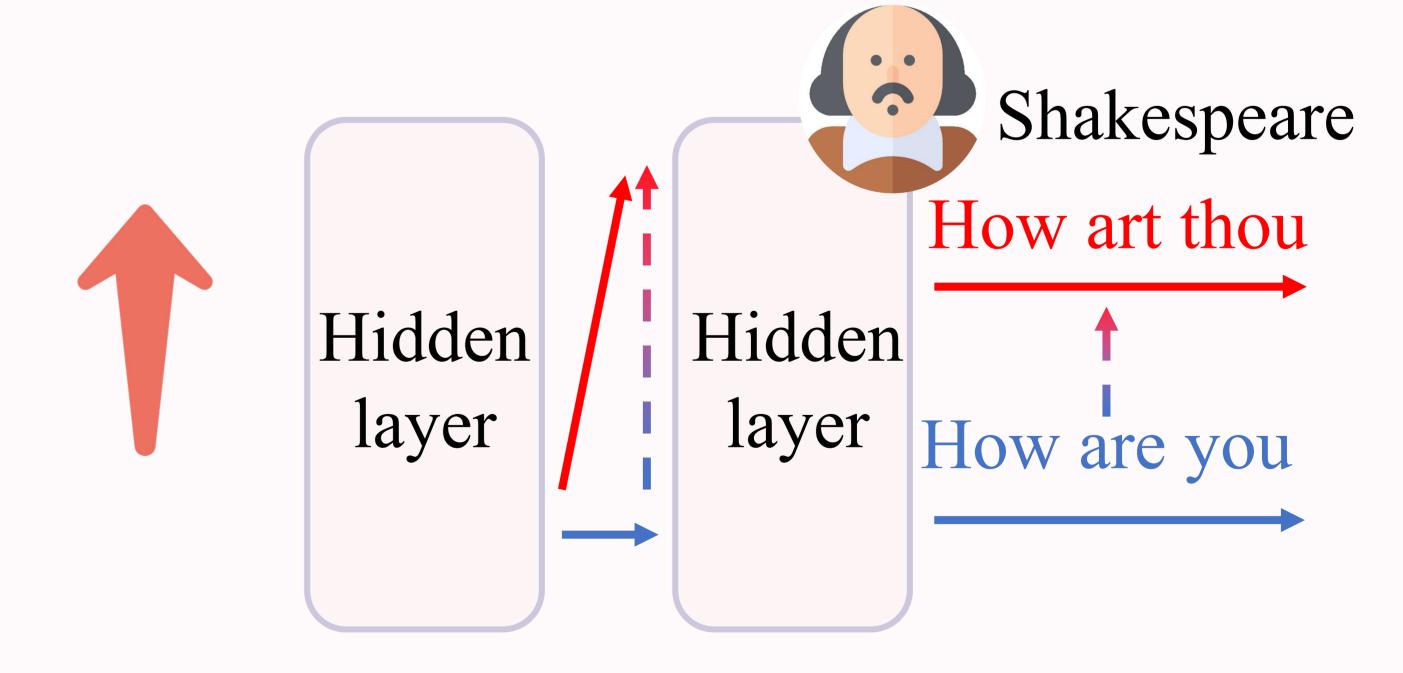


Inference Phase

6. Decoding-time Intervention



5. Latent Space Manipulation



4. Prompt Engineering