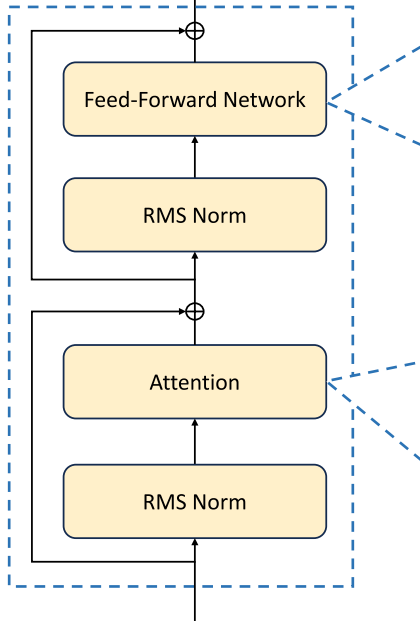
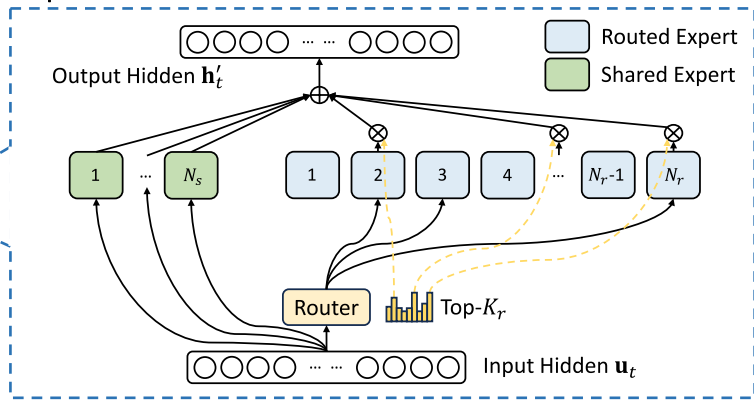


Transformer Block $\times L$



DeepSeekMoE



Multi-Head Latent Attention (MLA)

