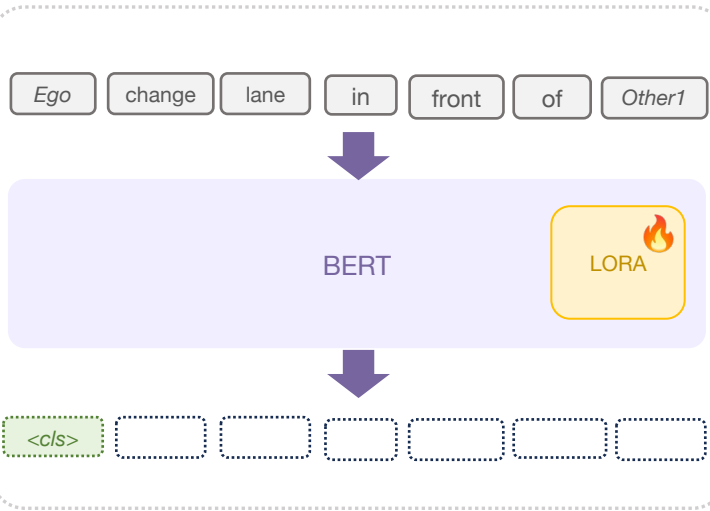
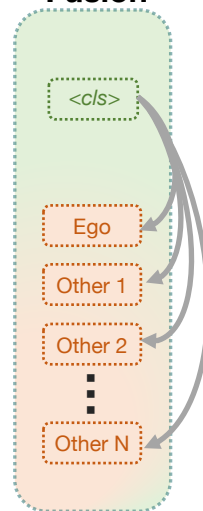


Language Encoder



Cross-Attn. Fusion



Text-Conditioned Diffusion Model

