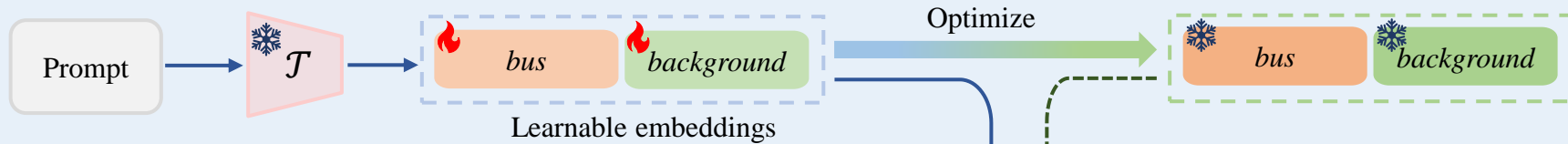
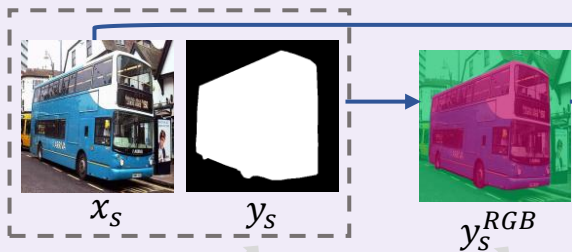


# Semantic Embedding Optimizer

Learnable Frozen

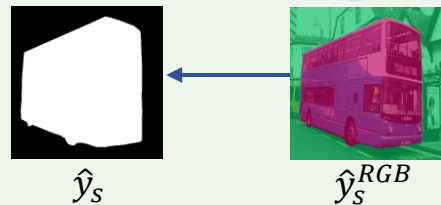


## Color Mask Transformer

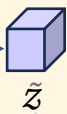
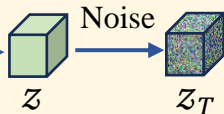


$\mathcal{L}_{\text{Mask}}$

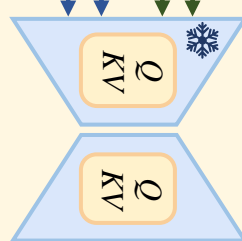
$\mathcal{L}_{\text{Pixel}}$



Noise



$\mathcal{L}_{\text{Latent}}$



Latent Diffusion Model



Binary Mask Extractor

$\rightarrow$  Support process  $\dashrightarrow$  Query process UNet Cross attention Image encoder Image decoder Text encoder