

Multimodal Facial Prompt Generator

A. Fine-grained Multimodal Feature Extractor

B. Overall Facial ID Feature Extractor

Enhanced Text

 Cross
Attention

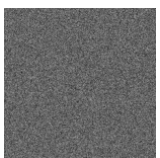
Image

 Cross
Attention

Denoising U-Net

$$\mathcal{L}_{total} = \mathcal{L}_{noise} + \mathcal{L}_{loc}$$

ID-Preservation network



+

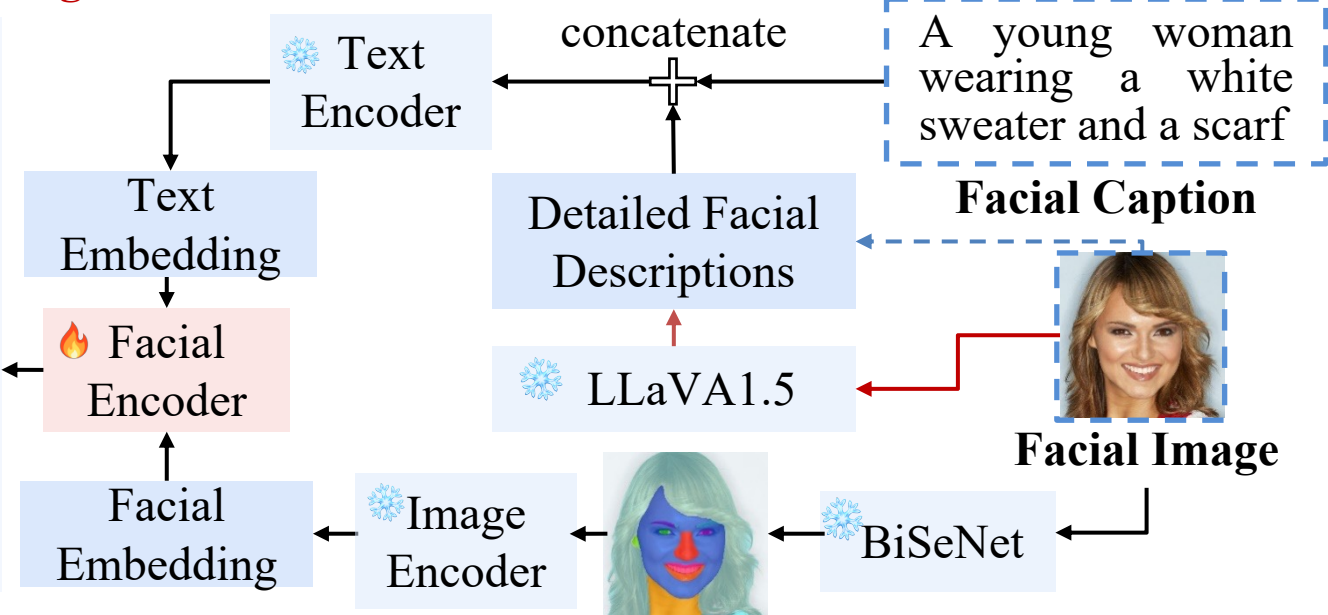


 Trainable Modules

 Frozen Modules

A. Fine-grained Multimodal Feature Extractor

Fine-grained Multimodal
Facial Feature



B. Overall Facial ID Feature Extractor

→ Training

- - - Inference

→ Shared

