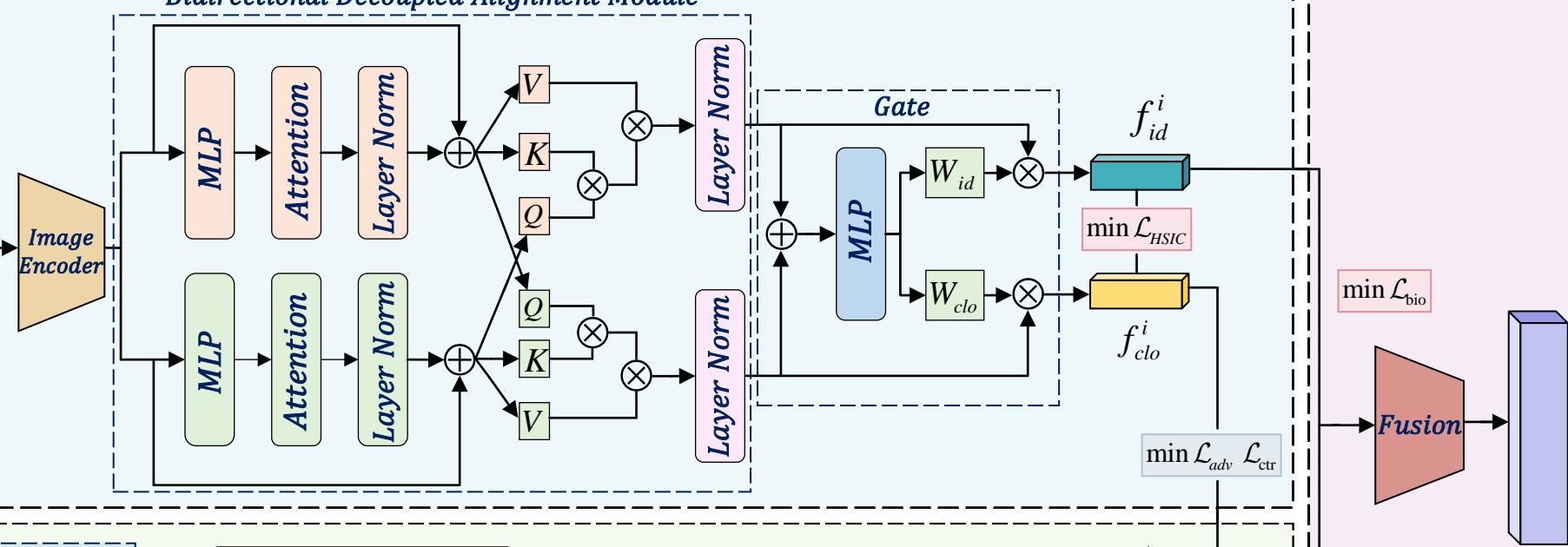


Bidirectional Decoupled Alignment Module



MLM

Clothes Caption:
"Dress, Hat, Shorts, Sleeve, Pants, Jacket, Coat, Skirt, Sweater, Sleeve, Cap, Scarf, Belt, Collar, Neck, Floral"
Person Caption:
"She wears a purple long sleeved, ankle length dress. There is a pattern on the dress."

Text Encoder

Attention

Layer Norm

Mean Pool

MLP

f_{clo}^t

f_{id}^t

$\min \mathcal{L}_{bio}$

