



Masked and Adaptive Transformer for Exemplar Based Image Translation

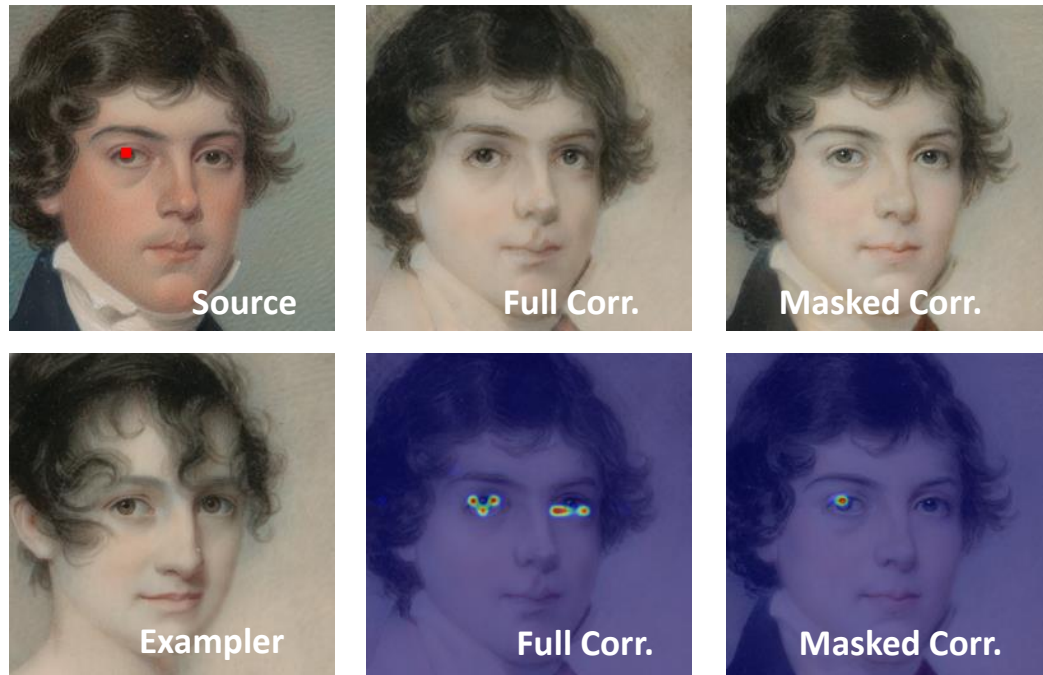
Chang Jiang¹, Fei Gao^{1,2*}, Biao Ma¹, Yuhao Lin¹, Nannan Wang², Gang Xu¹

¹Hangzhou Dianzi University ²Xidian University

<https://github.com/AiArt-HDU/MATEBIT>



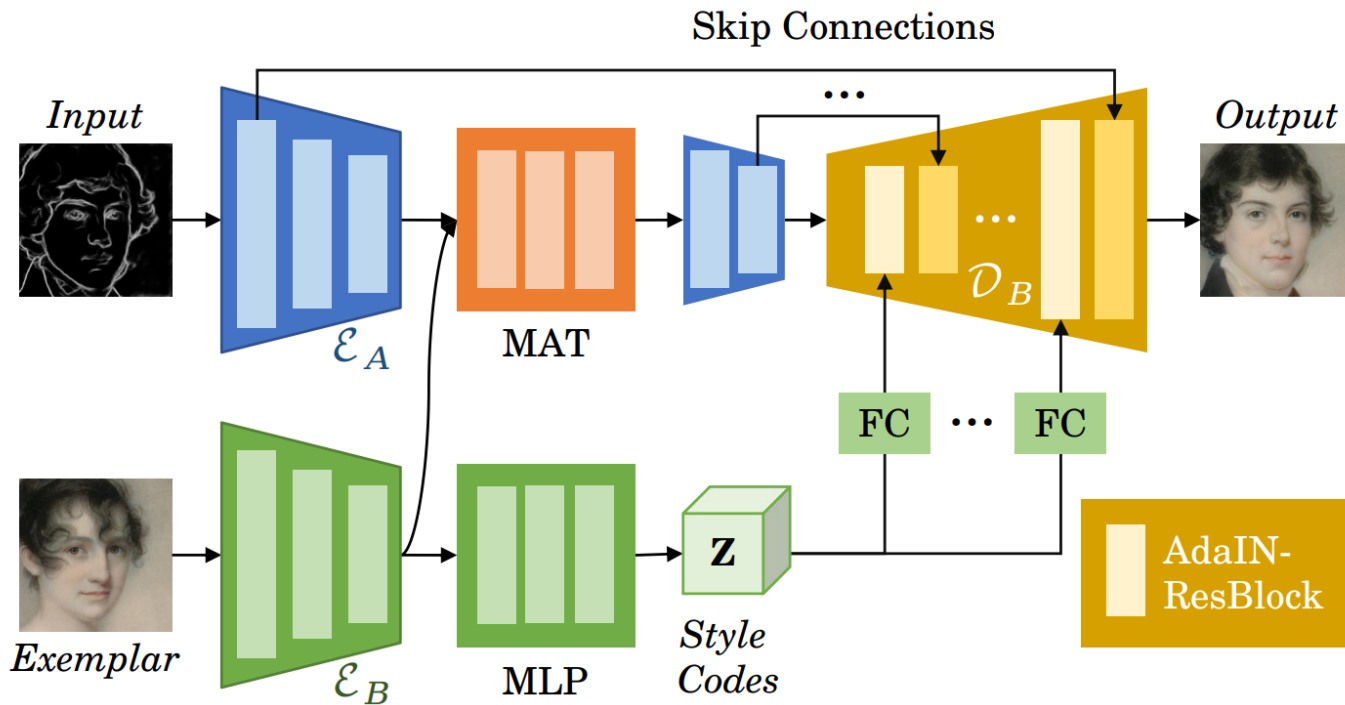
Motivation



- How to build accurate semantic correspondence between conditional inputs and exemplar ?
- How to modulate the style information for global injection?

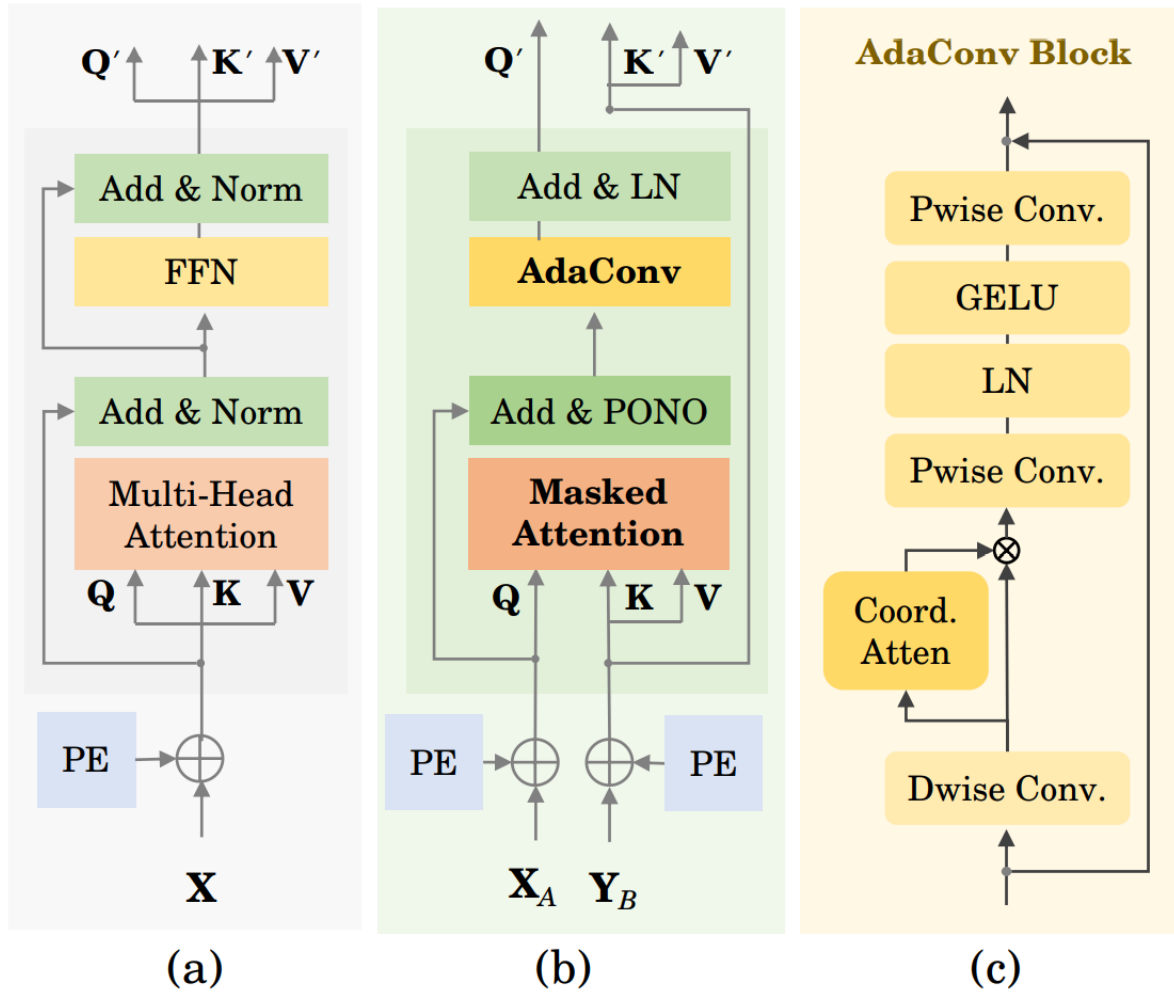
Improve the accuracy of matching on the one hand, and diminish the role of matching in image generation on the other hand.

Overview of translation network

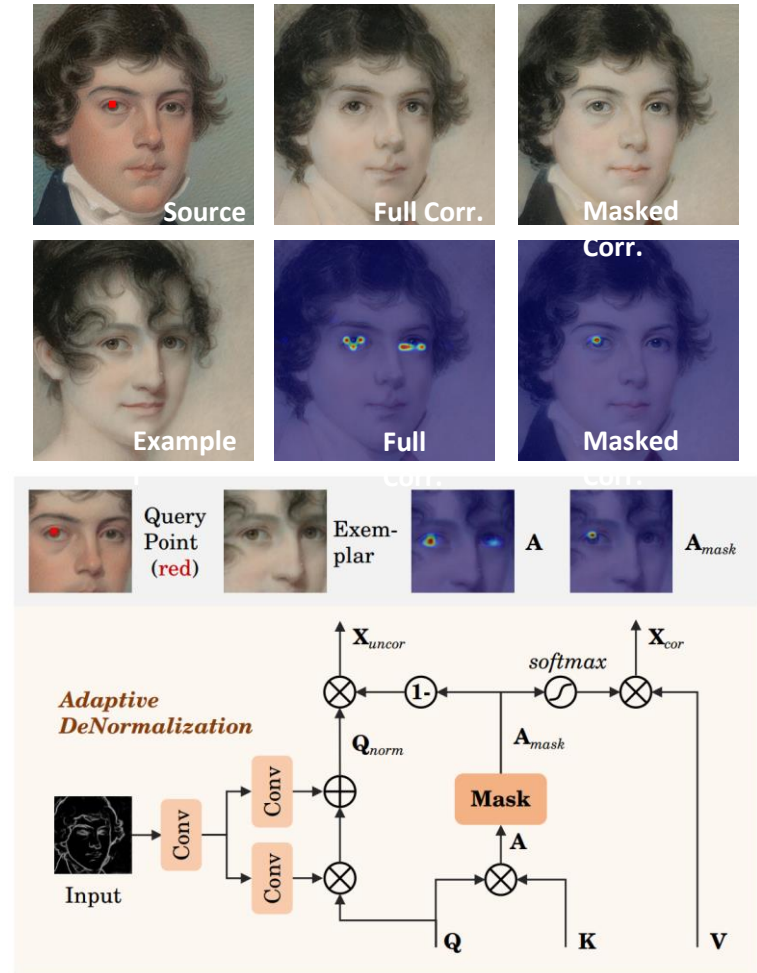


- Accurate semantic matching of correspondences
- Global style information modulation
- Significantly improve the image quality
- General translation solution

MAT Block



(a) Vanilla Transformer block, (b) MAT block, (c) AdaConv block



Masked Attention

MAT Block

Constructing conditional input and example image starting correspondence using the cosine attention mechanism.

$$A(u, v) = \frac{\tilde{Q}(u)\tilde{K}(v)^T}{\|\tilde{Q}(u)\| \|\tilde{K}(v)\|},$$

$$\tilde{Q}(u) = Q(u) - \bar{Q}(u), \tilde{K}(u) = K(u) - \bar{K}(u)$$

The mask interval is divided by the ReLU mechanism, A_{mask} is the post-mask probability matrix, and the post-mask confidence matrix is obtained by summing up the j dimensions.

$$X_{\text{cor}} = \tilde{A}_{\text{mask}}V, \text{ with } \tilde{A}_{\text{mask}} = \text{softmax}(\alpha \cdot A_{\text{mask}})$$

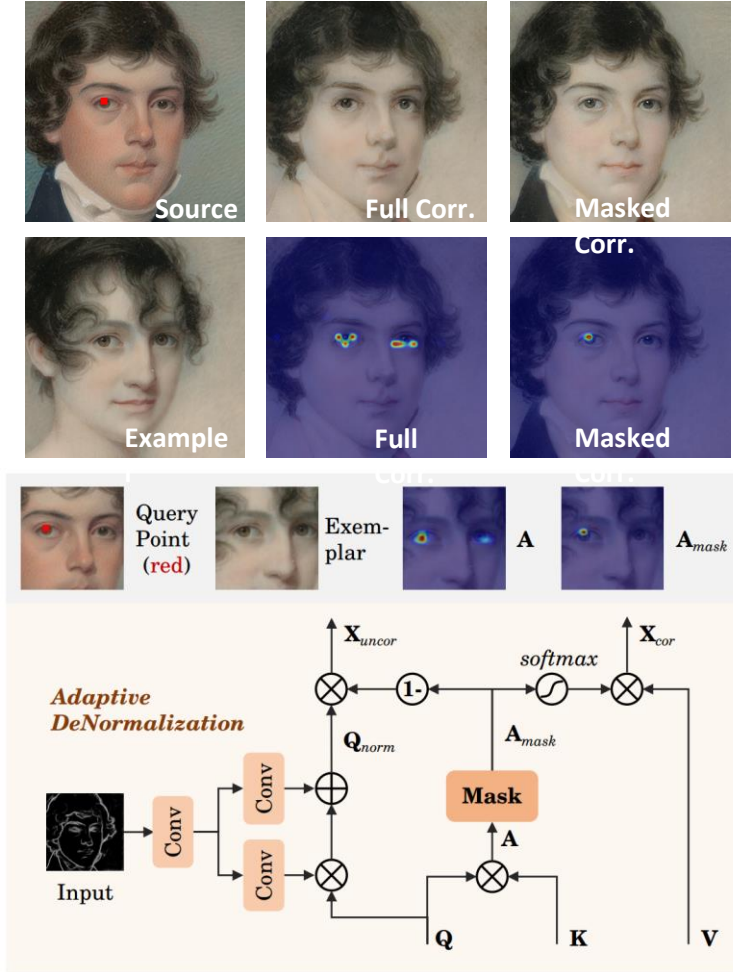
$$X_{\text{uncor}} = \left(1 - \sum_j A_{\text{mask}}\right) \odot Q_{\text{norm}}$$

$$Q_{\text{norm}} = \gamma(x_A) \frac{Q - \mu(Q)}{\sigma(Q)} + \beta(x_A)$$

$$X_{\text{agg}} = \text{PONO}(X_{\text{cor}} + X_{\text{uncor}} + Q)$$

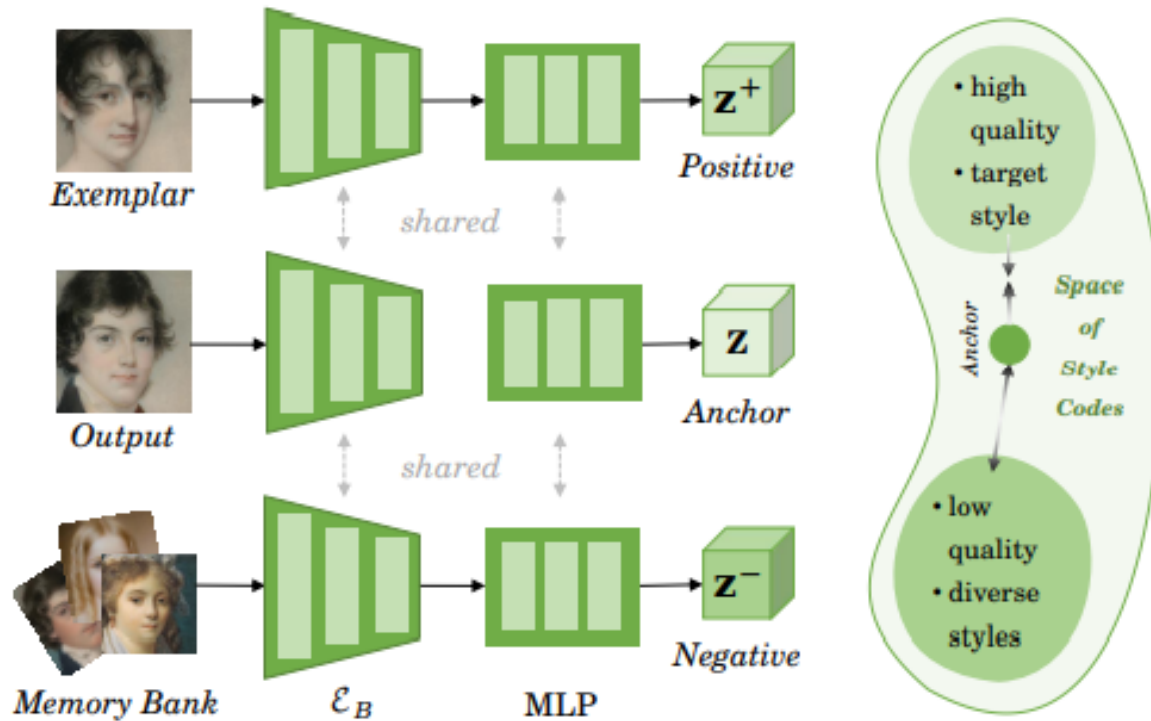
Fusing semantic information by adaptive convolution to generate a more detailed warp image.

$$X_{\text{MAT}} = \text{AdaConv}(X_{\text{agg}}) + X_{\text{agg}}$$



Masked Attention

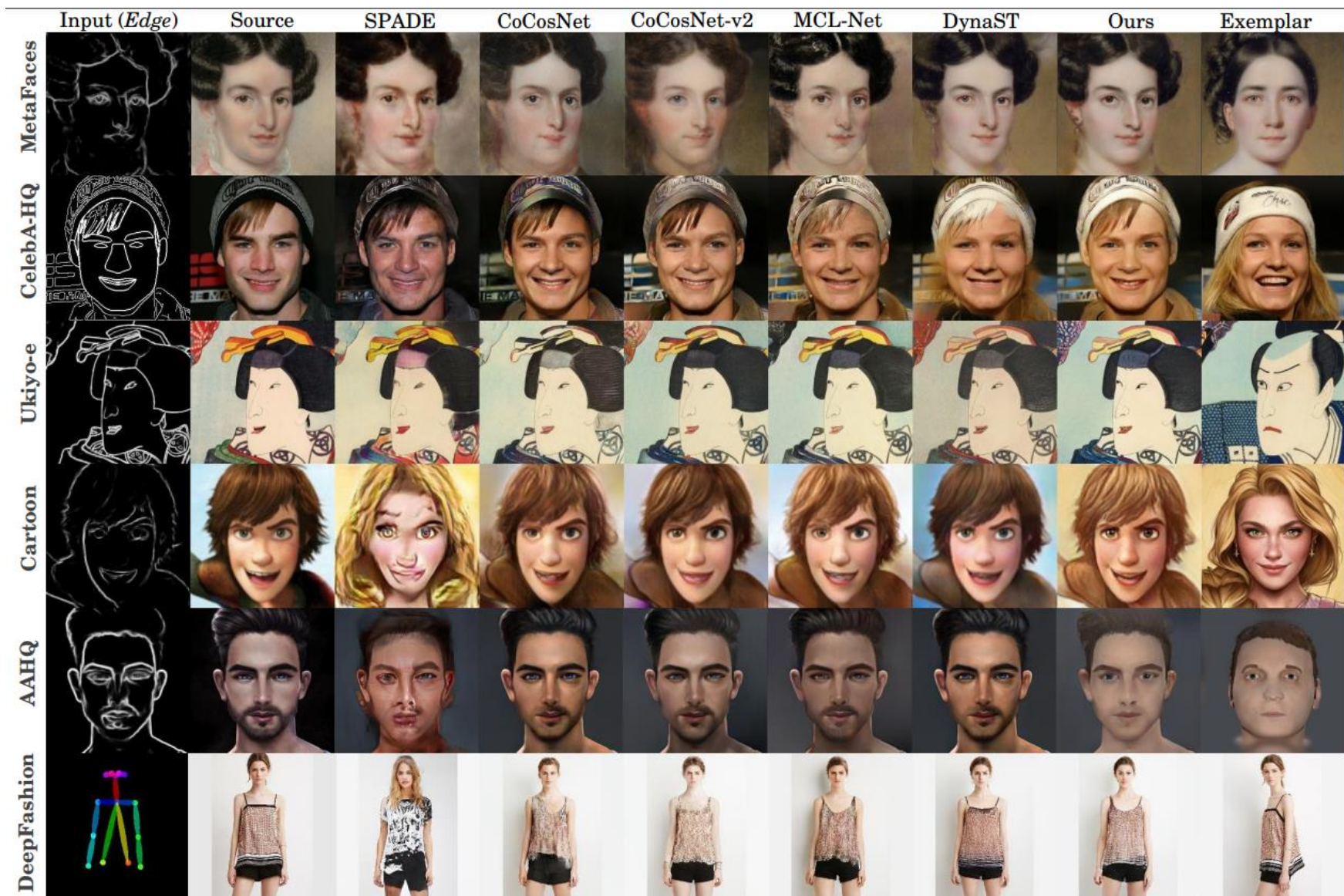
Contrastive style learning



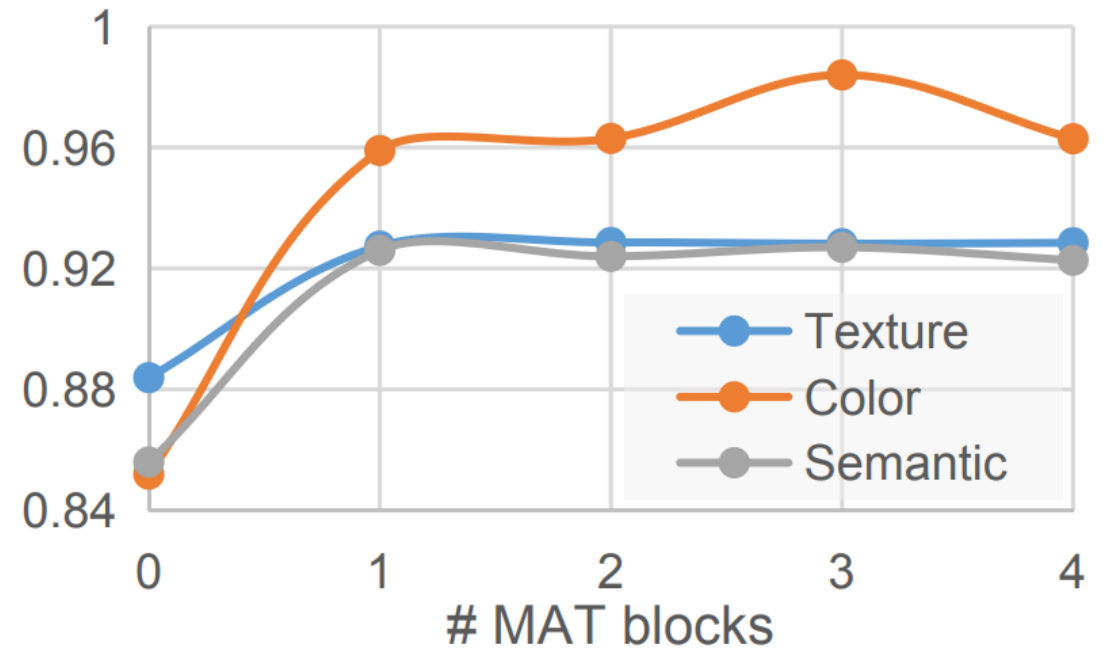
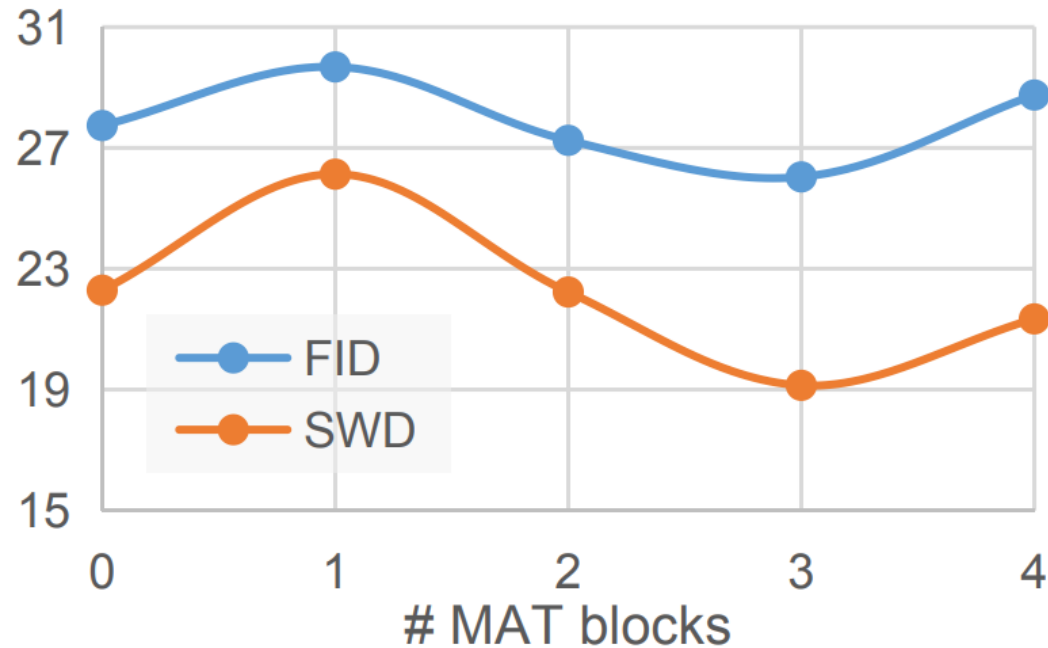
- Push to generate high-quality images
- Global style injection

$$\mathcal{L}_{style} = -\log \frac{\exp\left(\frac{\mathbf{z}^T \mathbf{z}^+}{\tau}\right)}{\exp\left(\frac{\mathbf{z}^T \mathbf{z}^+}{\tau}\right) + \sum_{j=1}^m \exp\left(\frac{\mathbf{z}^T \mathbf{z}_j^-}{\tau}\right)},$$

Comparison with state-of-the-art



Impacts of MAT



Impact of the number of MAT blocks on performance. For example, style, texture, semantic information, etc.

Conclusion: both the semantic consistency and style realism broadly improve with the number of MAT blocks and peak at three.

More Results



(a) Metfaces



(b) AAHQ



(c) CelebA-HQ



(d) Ukiyo-e-face

More Results



(e) Cartoon



(f) Chinese Ink Painting

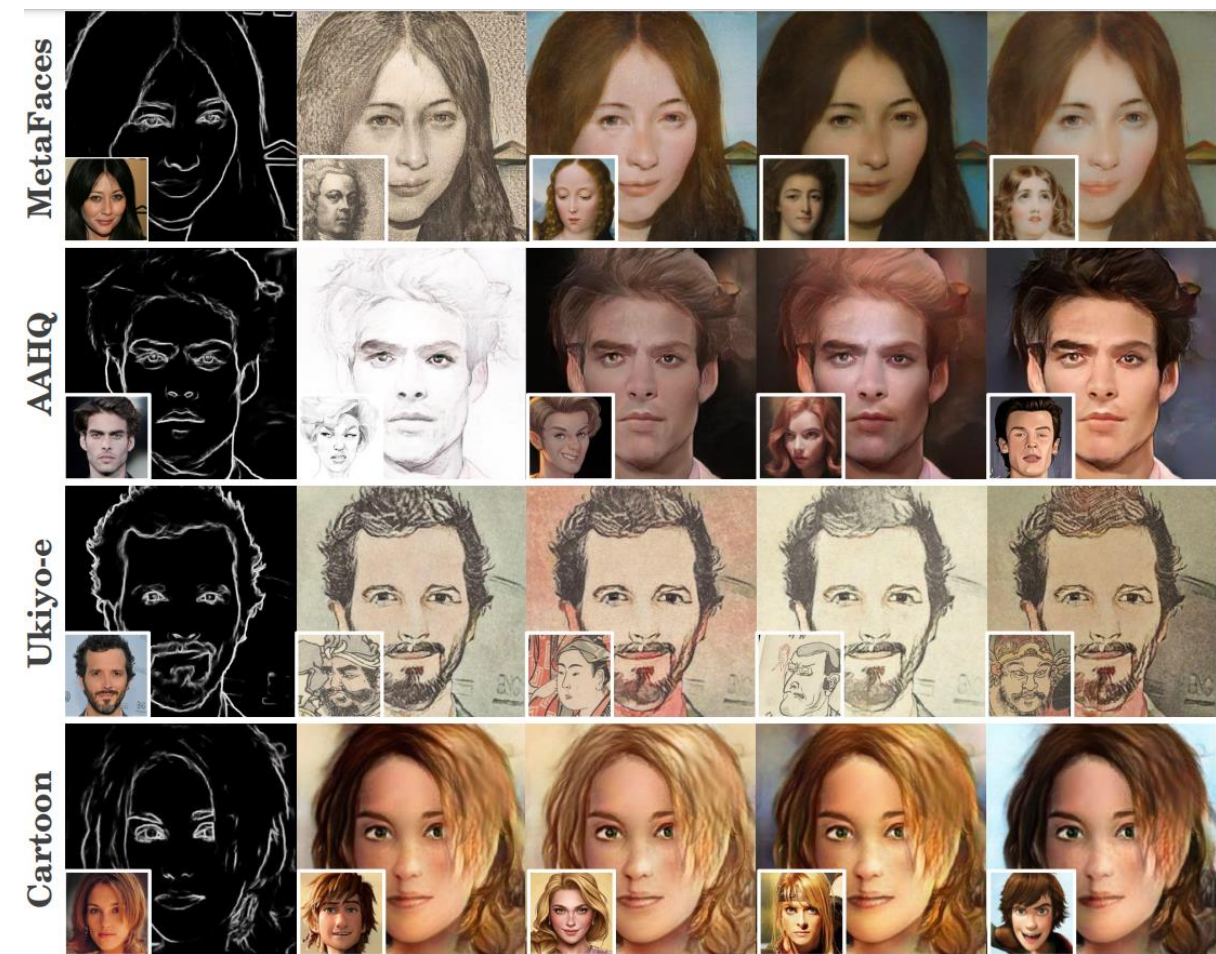


(g) DeepFashion



(h) Landscape

Applications



(a) Artistic Portrait Generation



Chinese ink paintings

(b) Chinese Ink Painting Generation