

# SLI-pSp: Injecting Multi-Scale Spatial Layout in pSp

Aradhya N. Mathur  
IIITD

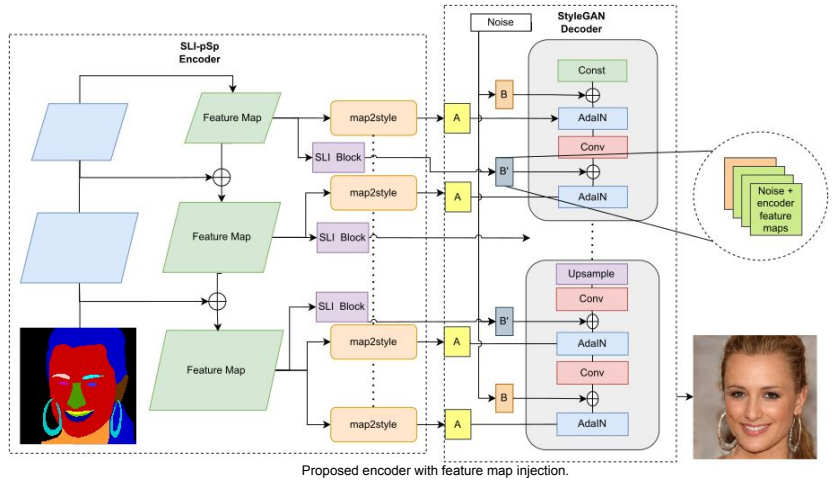
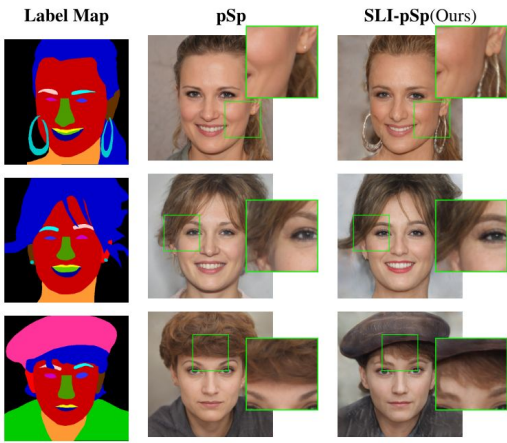
Anish Madan  
IIITD

Ojaswa Sharma  
IIITD

{aradhyam, anish16223, ojaswa}@iiitd.ac.in

## Abstract

We propose SLI-pSp, a general purpose Image-to-Image (I2I) translation model that encodes spatial layout information as well as style in the generator, using pSp as the base architecture. Previous methods like pSp have shown promising results by leveraging StyleGAN as a generator in various I2I tasks but they seem to miss finer or under-represented details in facial images like earrings and caps, and break down on complex datasets due to their solely global approach. We present a simple method to circumvent these issues without losing multi-modal properties of StyleGANs.



## Approach

- We propose a technique termed Spatial Layout Injection (SLI-pSp) that encodes spatial layout information in the input image in the StyleGAN generator along with style.
- We do so without modifying the style vector injection in the generator through pSp's *map2style* network, but rather by combining SLI with noise layers in the StyleGAN generator at multiple spatial scales.
- Such an approach helps preserve global aspects of image generation as well as enhance spatial layout details in the output. We experiment on several challenging datasets and across several I2I tasks that highlight the effectiveness of our approach over previous methods with respect to finer details in the generated image and overall visual quality.
- We extract feature maps using a feature pyramid network which are then propagated via *map2style* and SLI blocks and injected into the StyleGAN generator.

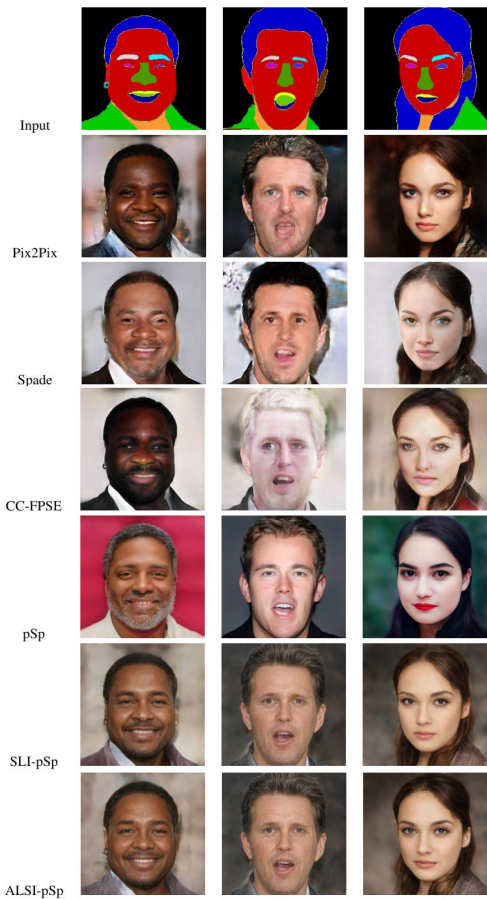
## Method

- Let the encoder feature maps be denoted as  $E_i$ , where  $i \in \{16, 32, 64\}$  corresponding to the spatial scales of the feature maps. These are generated using feature pyramid over a ResNet backbone.
- Let the noise layers in the StyleGAN generator be represented by  $N_j$ , where  $j \in \{4, 8, 16, \dots, 1024\}$  represents spatial sizes.
- The combined spatial layout feature maps and noise added to the generator,  $B'$  can be written as

$$B' = \text{concat}(\text{conv}(E_i), N_j) \quad i \in \{16, 32, 64\},$$

where *conv* represents a convolution layer and *concat* operation concatenates along the channel dimension as spatial sizes are the same.

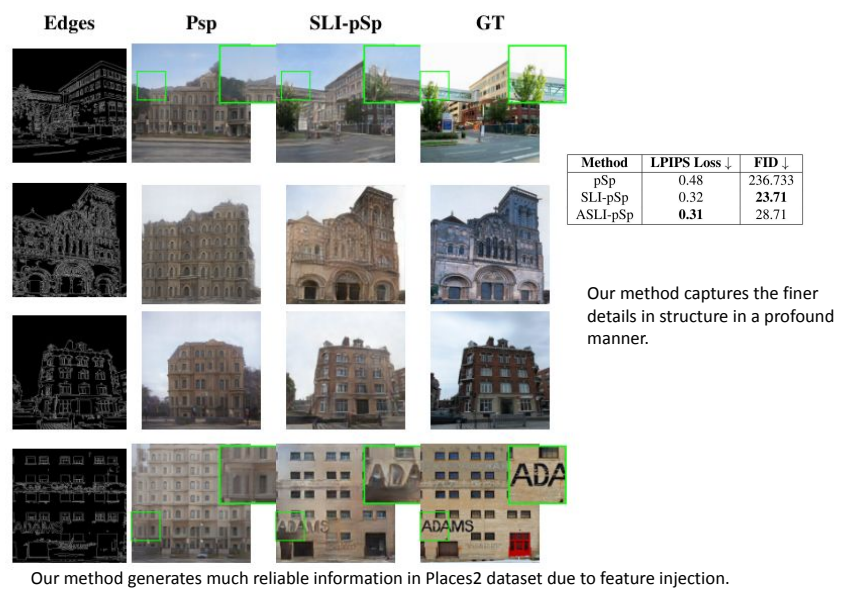
## Seg2Face



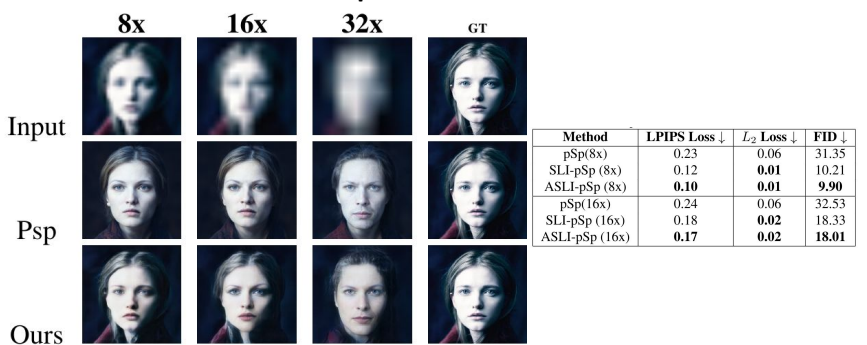
## Style Generation



## Edge2Image



## Super Resolution



## Conclusion

We propose a simple fix-ter-med Spatial Layout Injection (SLI) that encodes the spatial layout information from the input image and propagates it to the StyleGAN decoder with higher detail preservation.