# Learning to Generate Semantic Layouts for Higher Text-Image Correspondence in Text-to-Image Synthesis

**ICCV23** PARIS **KAIST**

Minho Park*   Jooyeol Yun*   Seunghwan Choi   Jaegul Choo
Korea Advanced Institute of Science and Technology (KAIST)

DAVIAN Data and Visual Analytics Lab

Project page

## Contribution

- We define a *Gaussian-categorical diffusion process* for modeling joint image-layout distributions, which is the first approach to unify two diffusion processes for image-layout generation.
- Our experiments reveal that *generating image-layout pairs can be a practical alternative to increase text-image correspondence* in circumstances where collecting web-scale text-image pairs is infeasible.
- We present *cross-modal outpainting*, which demonstrates that Gaussian-categorical diffusion models are also capable of modeling conditional distributions for semantic image synthesis and semantic segmentation.

## Motivation



Recall of facial attributes specified in the text descriptions.

*" Text-to-image generation approaches trained on small-scale dataset often fail to reflect text conditions. "*

## Gaussian-categorical Diffusion Process



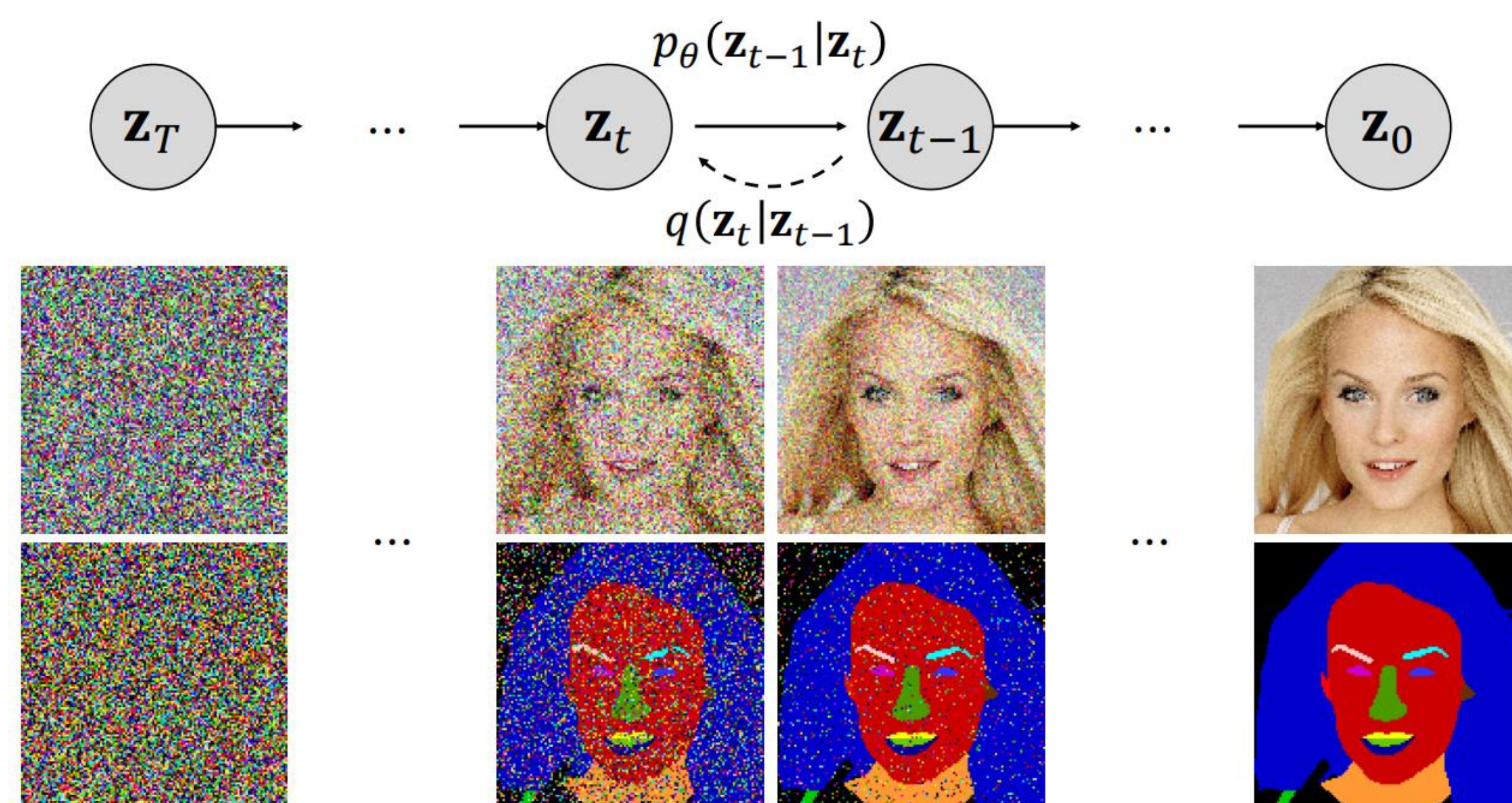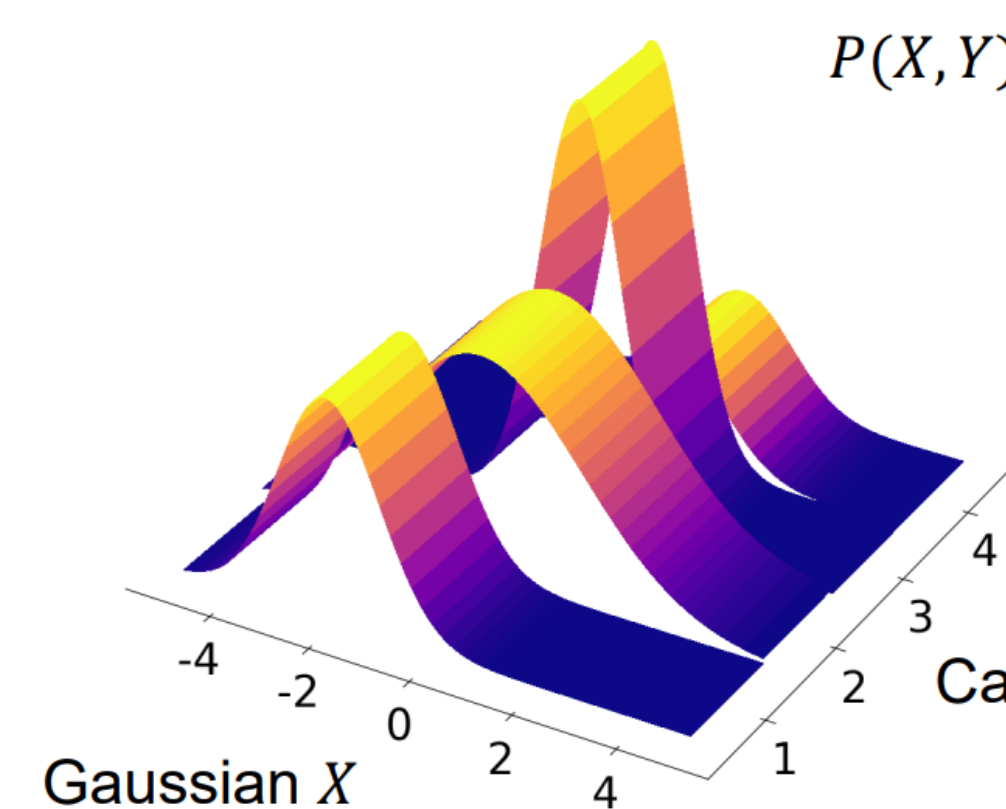*Illustration of the Gaussian-categorical diffusion process* on the image-layout distribution of MM CelebA-HQ.
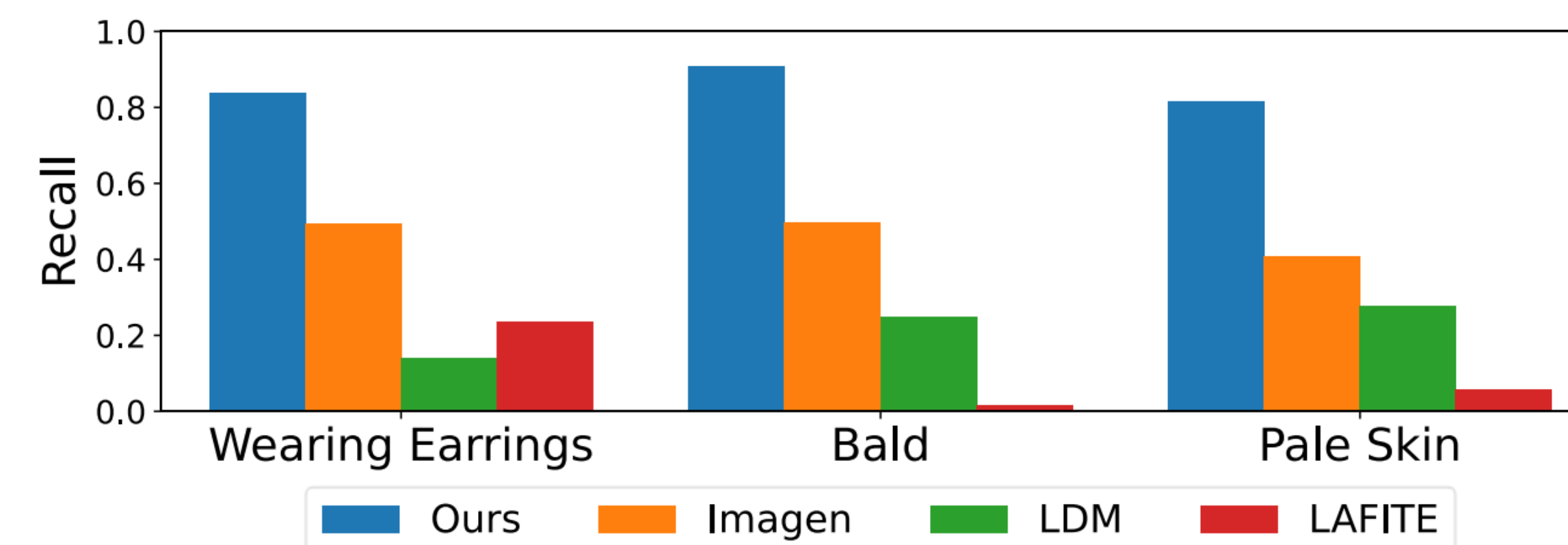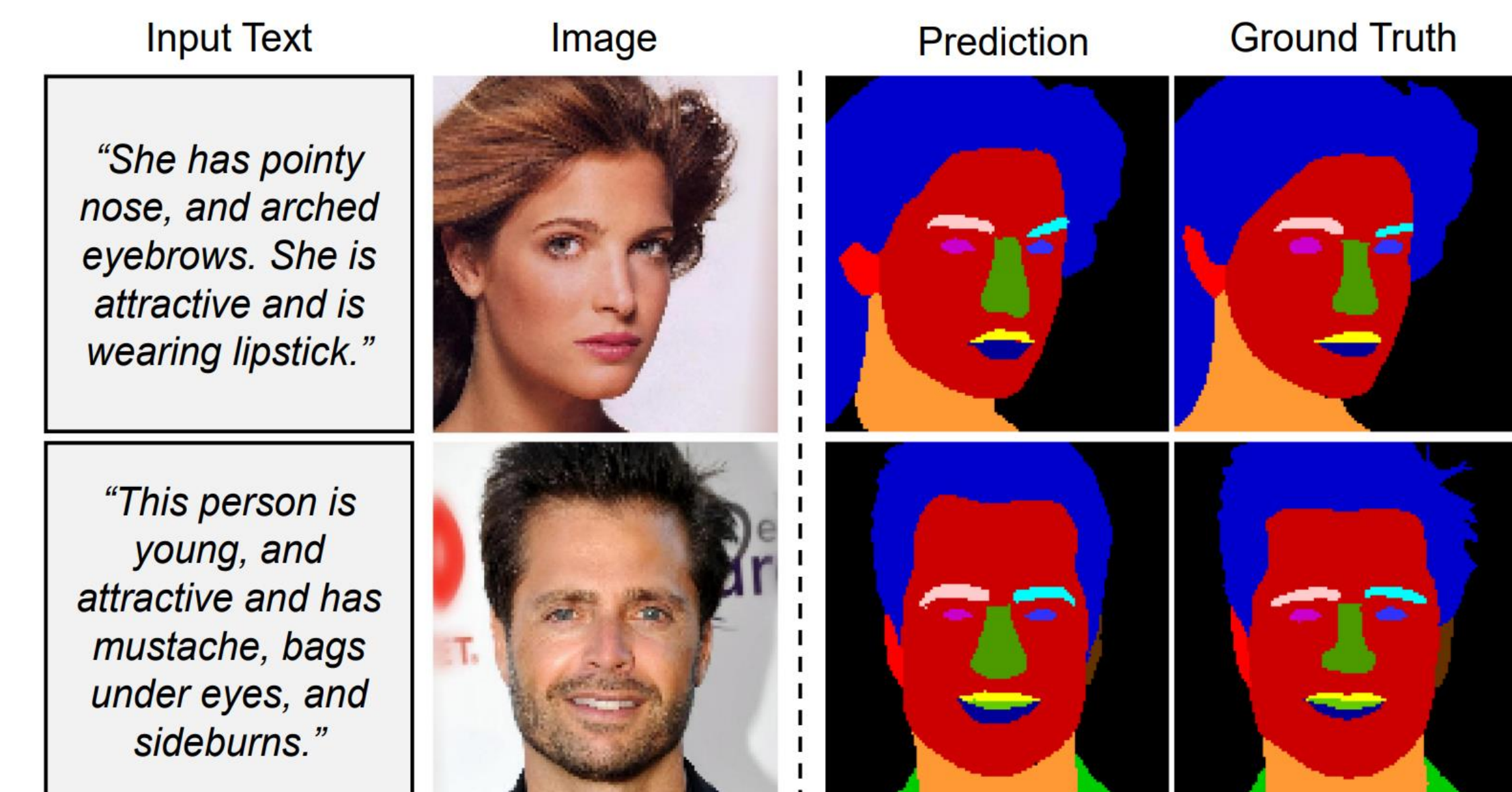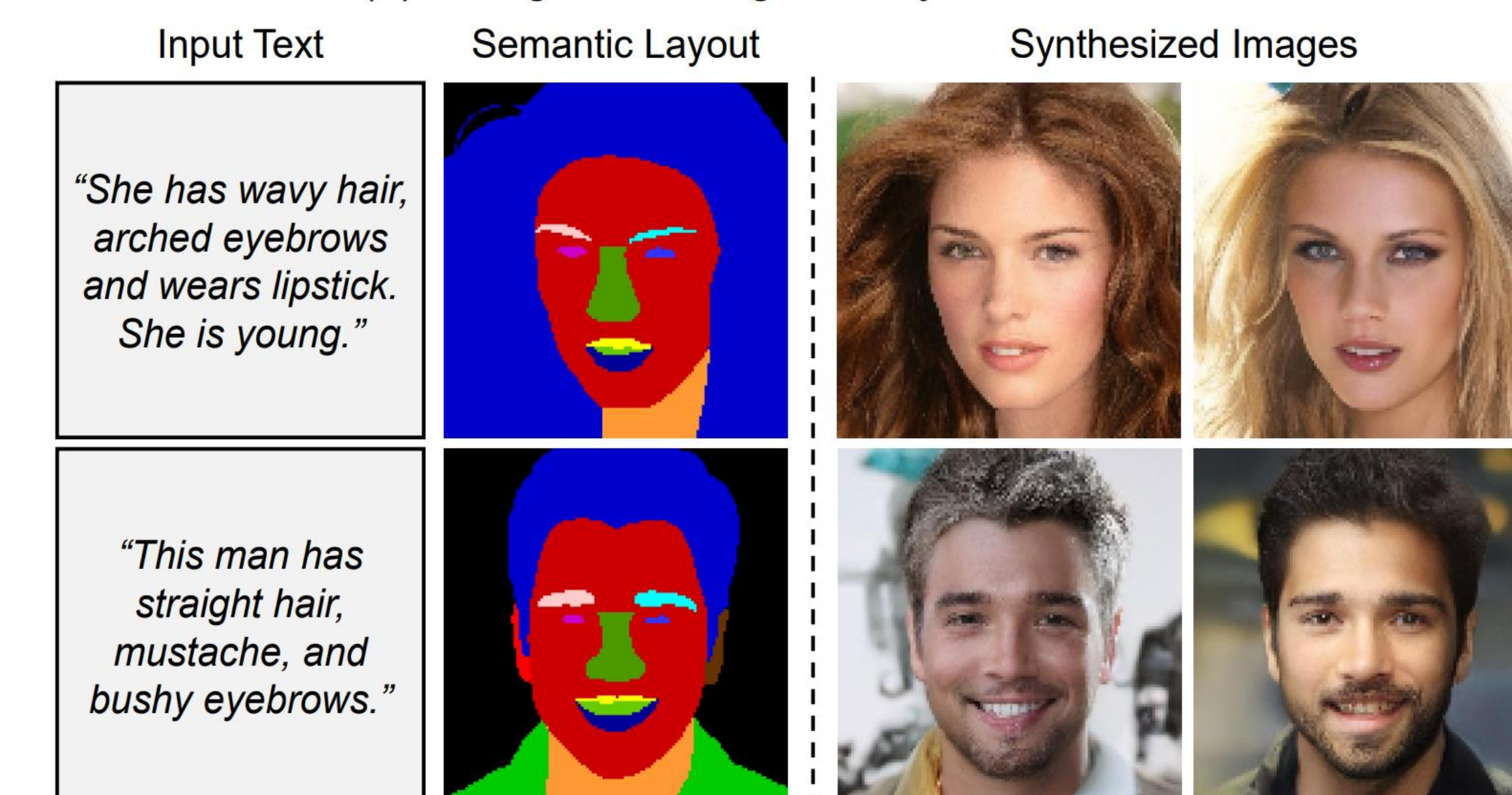
### Gaussian-categorical distribution



$$\mathcal{NC}(\mathrm{x},\mathrm{y};\mu,\Sigma,\Theta) = \mathcal{C}(\mathrm{y};\Theta)\cdot\mathcal{N}(\mathrm{x};\mu_{\mathrm{y}},\Sigma_{\mathrm{y}})$$

where $\mathrm{x}\in\mathbb{R}^N, \mathrm{y}\in\{1,2,\ldots,K\}^M\subset\mathbb{R}^M$
$\mu\in\mathbb{R}^{S\times N}, \Sigma\in\mathbb{R}^{S\times N\times N}, \Theta\in\mathbb{R}^{M\times K}, (S=K^M)$
$\mu_{\mathrm{y}}\in\mathbb{R}^N, \Sigma_{\mathrm{y}}\in\mathbb{R}^{N\times N}$

### Final objective function

$$D_{KL}\big(q(\mathbf{z}_{t-1}|\mathbf{z}_t,\mathbf{z}_0)\parallel p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t)\big) = \mathbb{E}_q\left[\frac{1}{2\sigma_t^2}\|\widetilde{\boldsymbol{\mu}}_t - \widetilde{\boldsymbol{\mu}}_\theta(\mathbf{z}_t)\|^2\right] + D_{KL}\big(\widetilde{\Theta}_t\parallel\widetilde{\Theta}_\theta(\mathbf{z}_t)\big) + C$$

Gaussian Diffusion     Categorical Diffusion

* Detailed proofs for each step are provided in A.1.
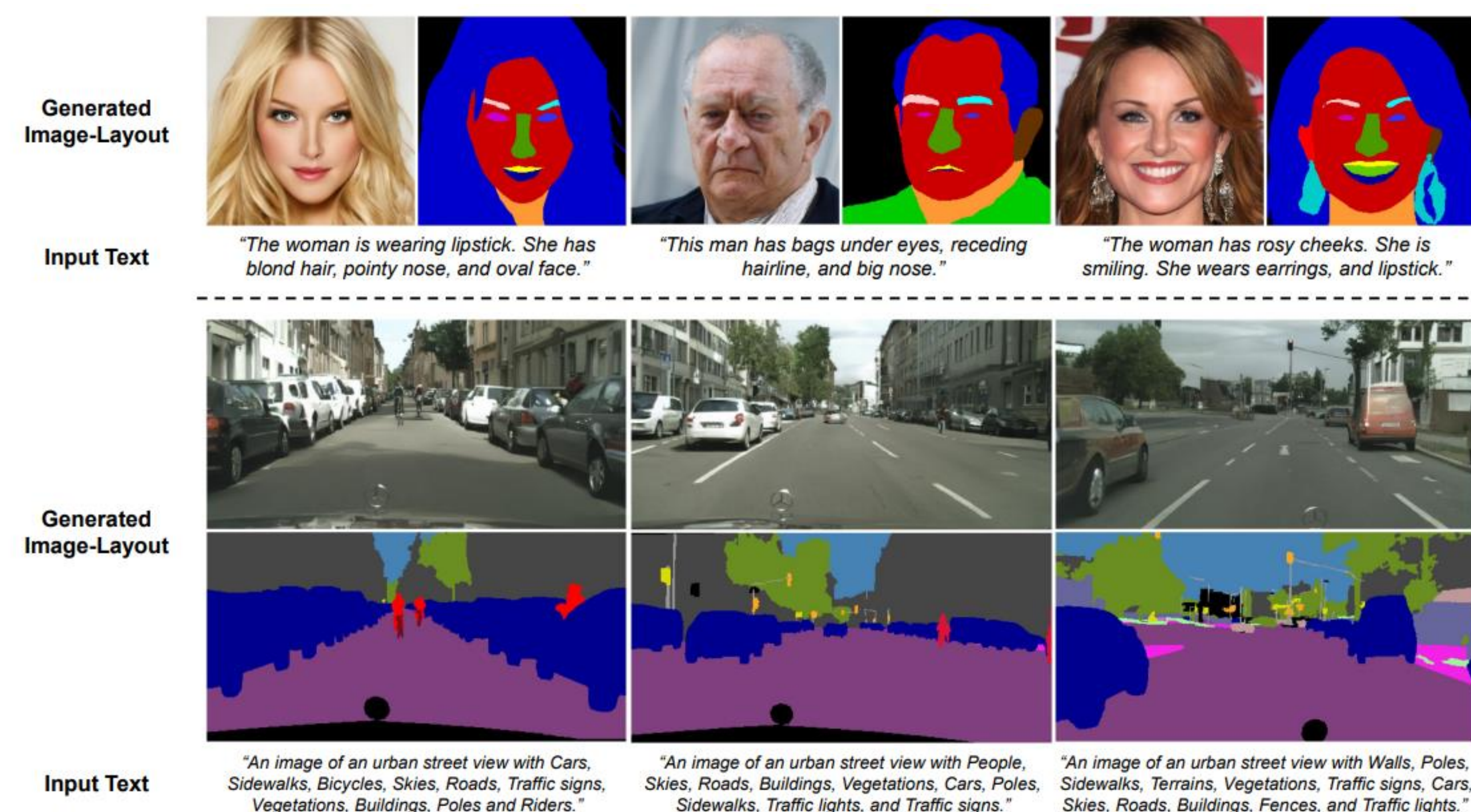
## Cross-modal Outpainting



(a) Text-guided Image-to-Layout Generation
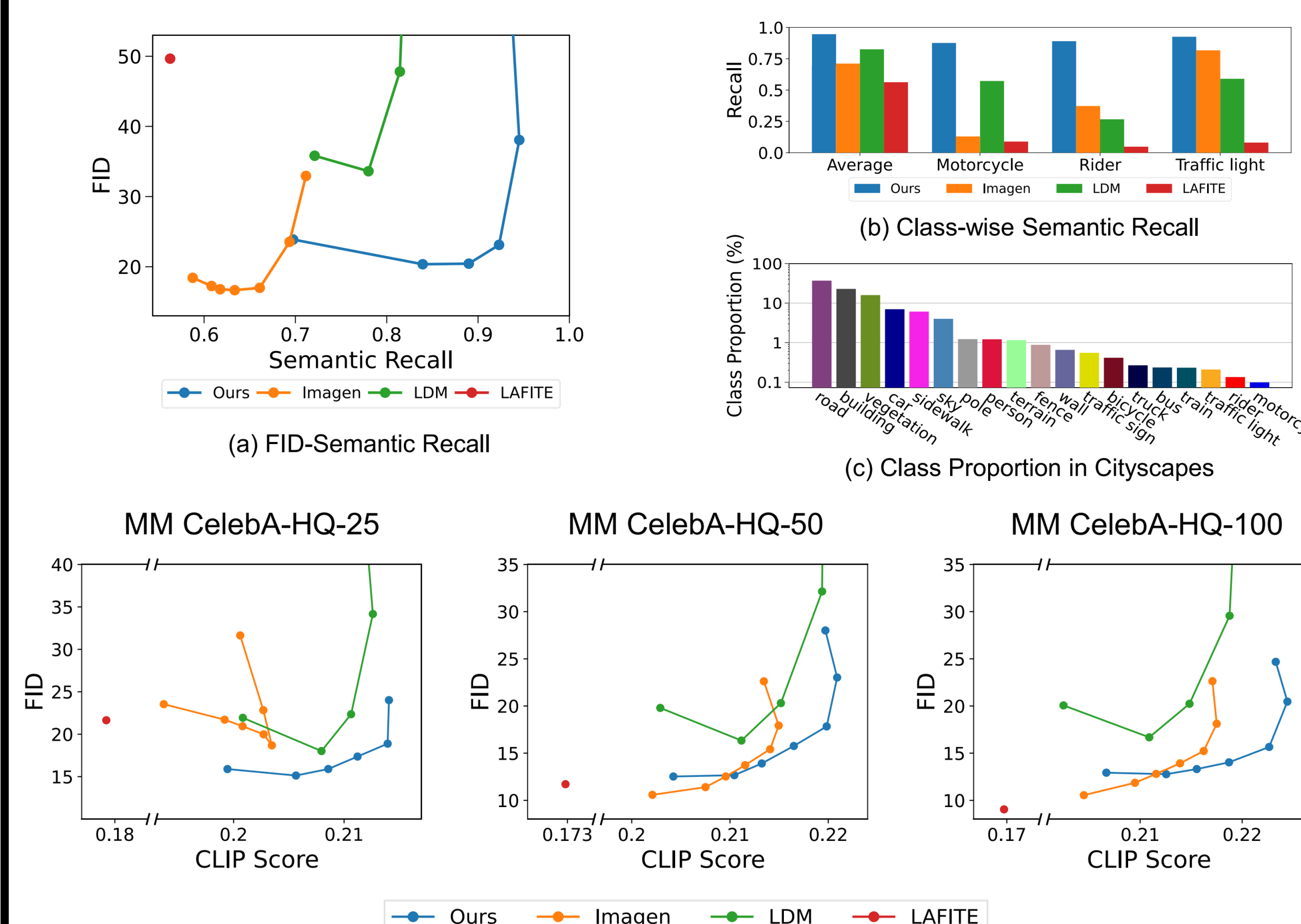


(b) Text-guided Layout-to-Image Generation

*" Since we model image-layout joint distribution, we can also model conditional distributions. "*
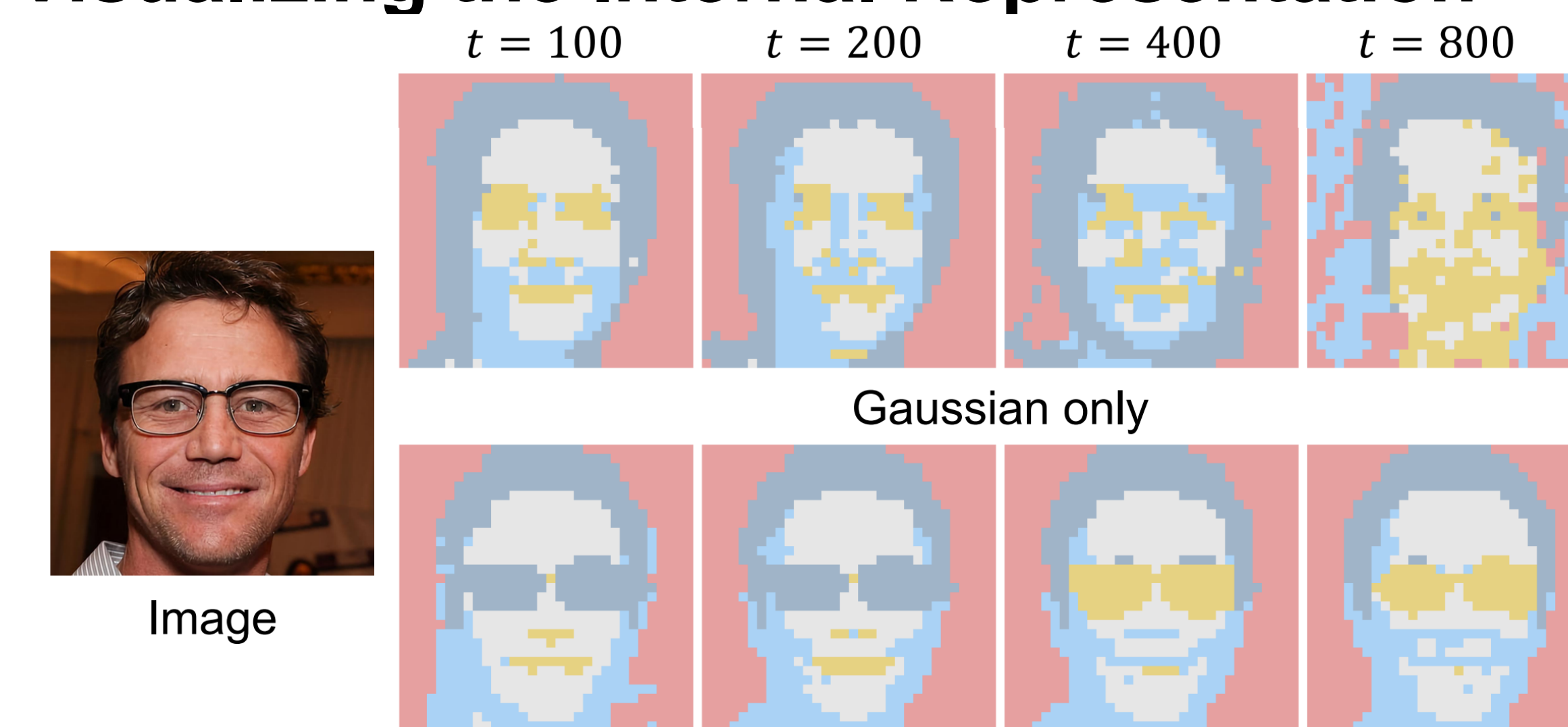
## Qualitative Results



*" GCDP generates aligned image-layout pairs from text descriptions. "*

## Quantitative Comparison



(a) FID-Semantic Recall
(b) Class-wise Semantic Recall
(c) Class Proportion in Cityscapes

MM CelebA-HQ-25     MM CelebA-HQ-50     MM CelebA-HQ-100

*" GCDP increase text-image correspondence maintaining image quality. "*

## Visualizing the Internal Representation



Gaussian only

Gaussian-categorical

*" GCDP has better understanding of semantic layout than Gaussian-only Diffusion Model. "*

## Find us!

Jooyeol Yun     Minho Park