



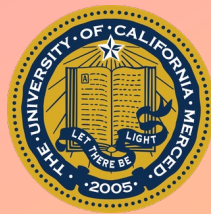
Omni-Attribute:

Open-vocabulary Attribute Encoder for Visual Concept Personalization

Tsai-Shien Chen^{1,2} Aliaksandr Siarohin¹ Gordon Guocheng Qian¹ Kuan-Chieh Jackson Wang¹ Egor Nemchinov¹
Moayed Haji-Ali¹ Riza Alp Guler¹ Willi Menapace¹ Ivan Skorokhodov¹ Anil Kag¹ Jun-Yan Zhu³ Sergey Tulyakov¹



¹Snap Inc.



²UC Merced



³CMU



Copy-Paste Artifacts

Phantom [ICCV'25]



Movie Weaver [CVPR'25]



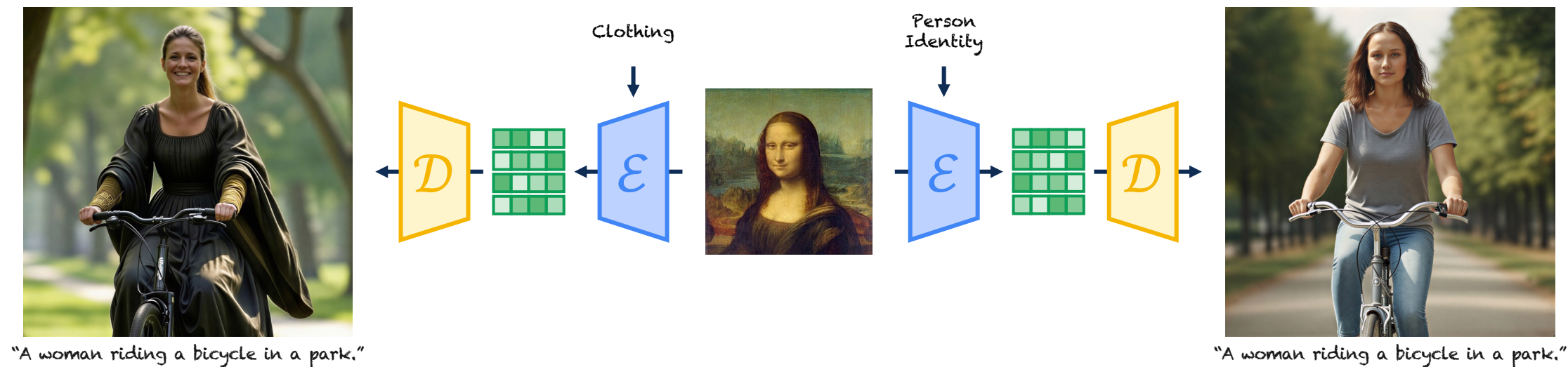
BindWeave [arXiv'25]



Can we prevent such information leakage at the encoder stage?

Can we learn an attribute encoder that only extracts the information of the user-specified attributes?

Open-vocabulary Image Attribute Encoder



Interactive Demo



Try it yourself on the [project webpage!](#)

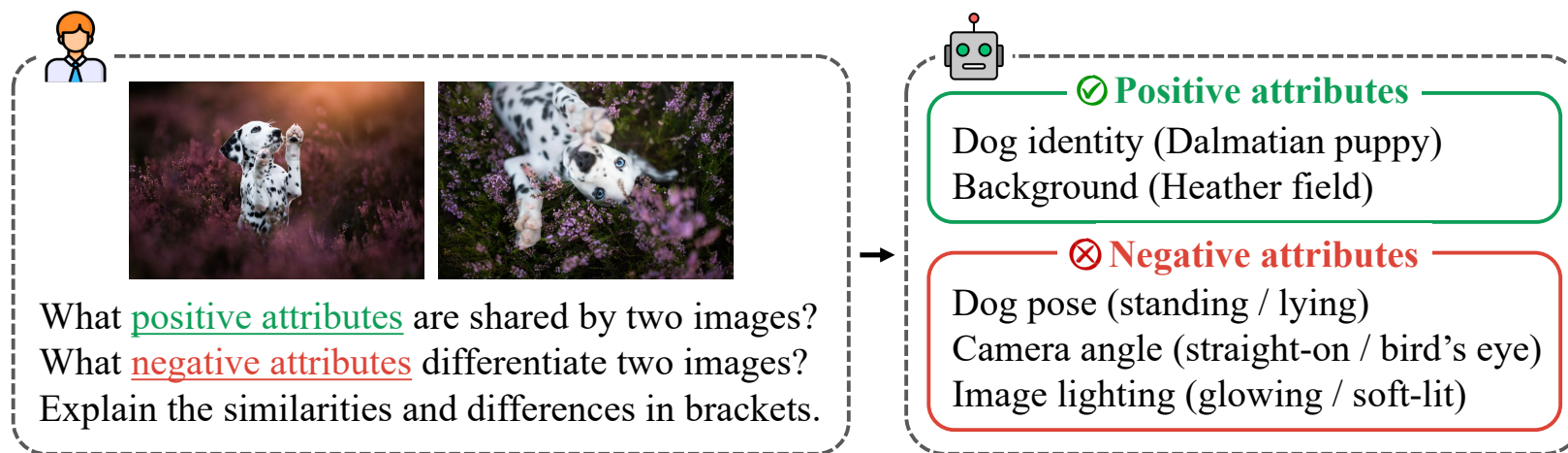


Attribute-level Image Composition



Person Identity	Person Identity	Clothing	Person	Hat	Lighting	Expression	Hairstyle	Makeup
"A person riding a roller coaster."		"A man in a shirt sitting at a bar."		"A man with a hat looking at the camera."		"A close-up view of a woman making a face."		
Artistic Style	Knight	Pose	Dog	Clothing	Tone	Rabbit	Background	Art Style
"A person riding a roller coaster."		"A knight posing against a gray backdrop."		"A dog in a kimono walking on the sand."		"A rabbit lying in the snow."		

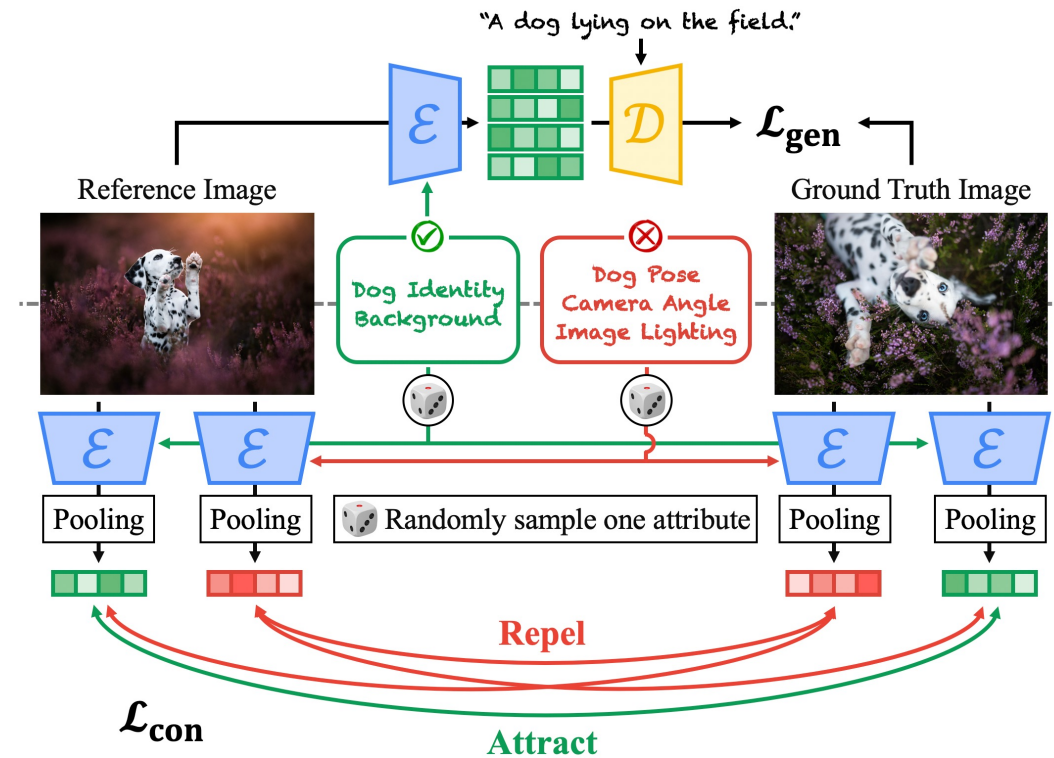
Semantic Connections between Image Pairs



Attribute-level Representation Learning

[Goal 1] The embeddings need to capture sufficient information to ensure high-fidelity reconstruction of the target attributes.

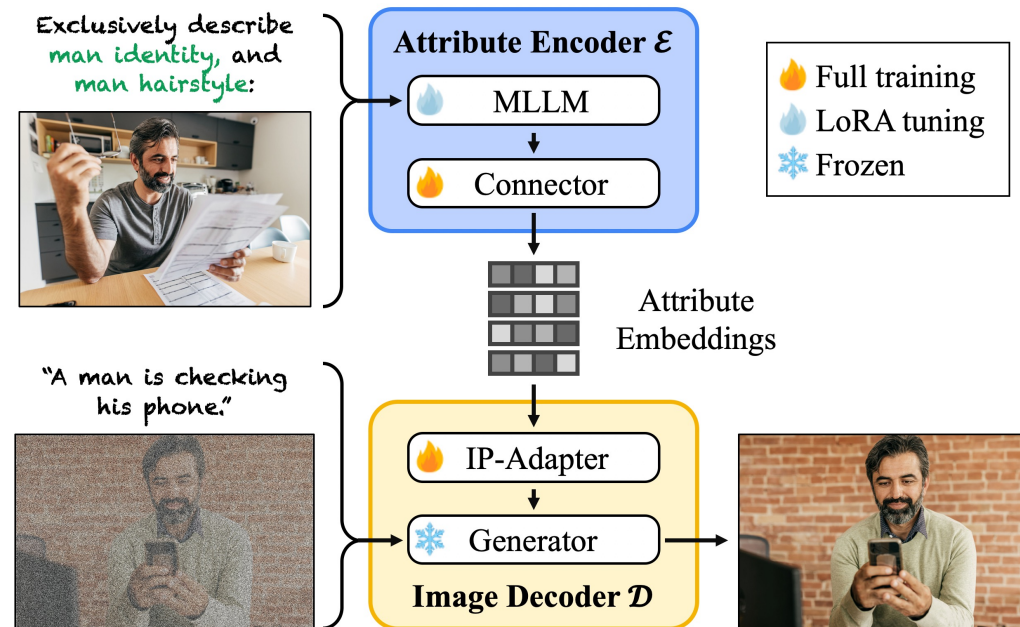
[Goal 2] The embeddings need to remove irrelevant information from other attributes.



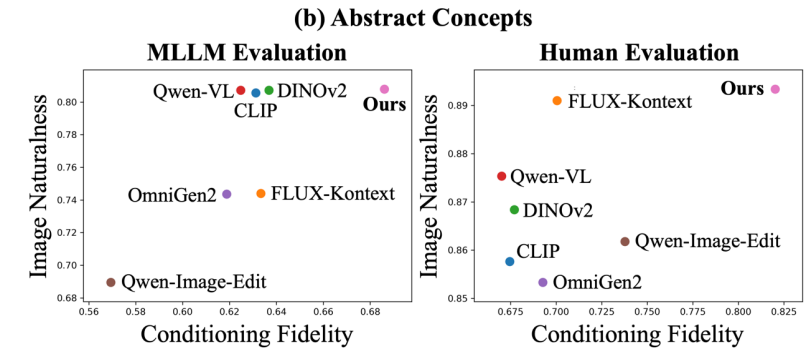
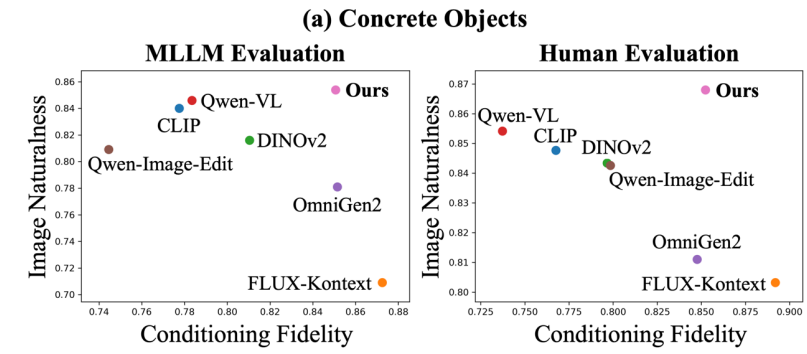
MLLM-based Attribute Encoder

[Goal 1] Capability to jointly process both image and text inputs.

[Goal 2] Strong vision-language prior to support attribute disentanglement.



Open-vocabulary Attribute Personalization



*For image encoders, we train IP-Adapter modules between each encoder and the same frozen image generation backbone.

*For image editing models, we reformulate the prompt as "Preserve the <attribute> of the image and generate <prompt>" to enable attribute personalization.

Composability of Attribute Embeddings

Motivated by Composable Diffusion [ECCV'22]:

$$\hat{\epsilon}(\mathbf{x}_t, t) = \underbrace{\epsilon_\theta(\mathbf{x}_t, t)}_{\text{Unconditional prediction}} + \sum_{i=1}^n w_i \underbrace{(\epsilon_\theta(\mathbf{x}_t, t | \mathbf{c}_i) - \epsilon_\theta(\mathbf{x}_t, t))}_{\text{Conditional score direction (cond pred - uncond pred)}}$$

Adaption to Flow-matching Model:

Conditional flow field of **single** image-attribute condition: $\Delta_{(I_i, a_i)} = \mathcal{D}(\mathcal{E}(I_i, a_i), \emptyset) - \mathcal{D}(\emptyset, \emptyset)$

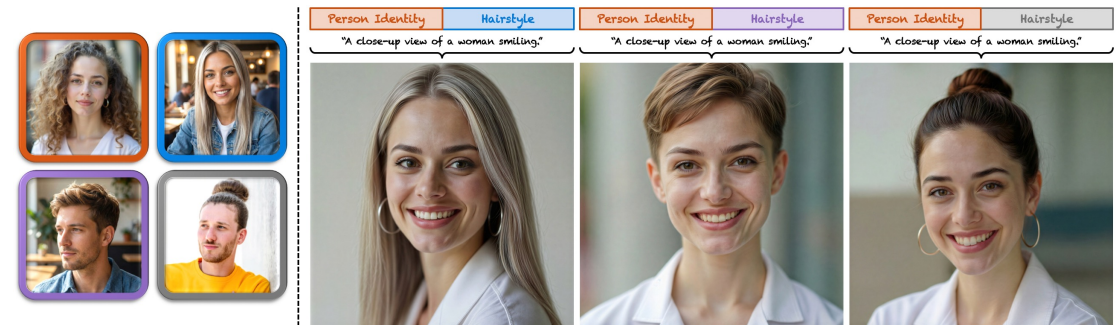
Composition of **multiple** image-attribute conditions: $v^* = \mathcal{D}(\emptyset, c) + \sum_{i=1}^N w_i \cdot \Delta_{(I_i, a_i)}$

Practical Applications

Advertisement Image Synthesis



Hairstyle Customization



Storytelling Visualization



Creative Content Generation



Omni-Attribute:

Open-vocabulary Attribute Encoder for Visual Concept Personalization

Project Page

