

SEGA:

Instructing Text-to-Image Models using Semantic Guidance

Manuel
BrackFelix
FriedrichDominik
HintersdorfLukas
StruppekPatrick
SchramowskiKristian
Kersting

Semantic Control over Diffusion Models

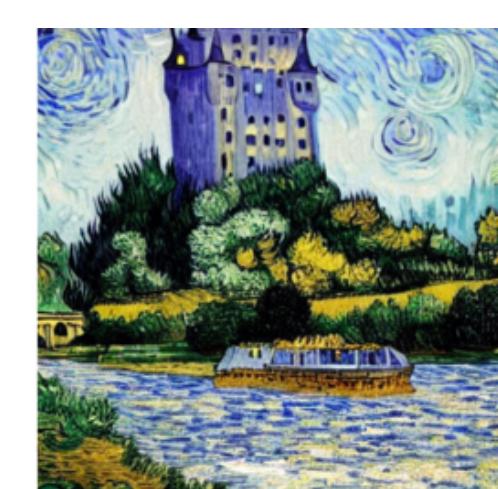
- Producing **high-fidelity** images with text-to-image models is an **iterative, interactive** process
- Users want to make **changes** until **satisfied** with the outputs
- Fine-grained **semantic control** is needed which should be as **easy to use** and **versatile**

Method	Manipulation Success
Disentanglement (Wu, 2022)	41.35 %
Prompt2Prompt (Hertz, 2023)	43.25 %
Composable Diffusion (Liu, 2022)	60.50 %
SEGA (Ours)	72.72 %

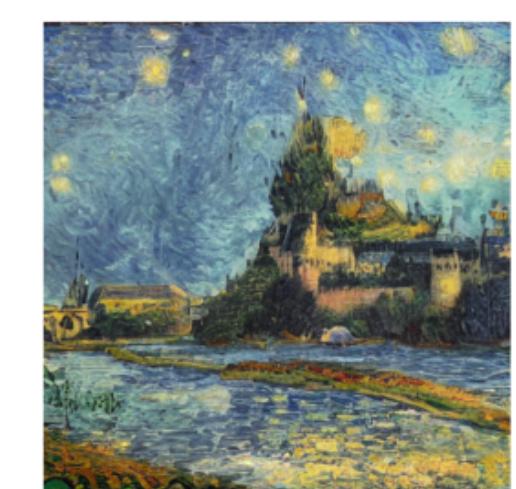
User Preference

Semantic
Guidance

83% vs 13%



Original
+ 'van Gogh's Starry Night'
+ 'boat in the river'

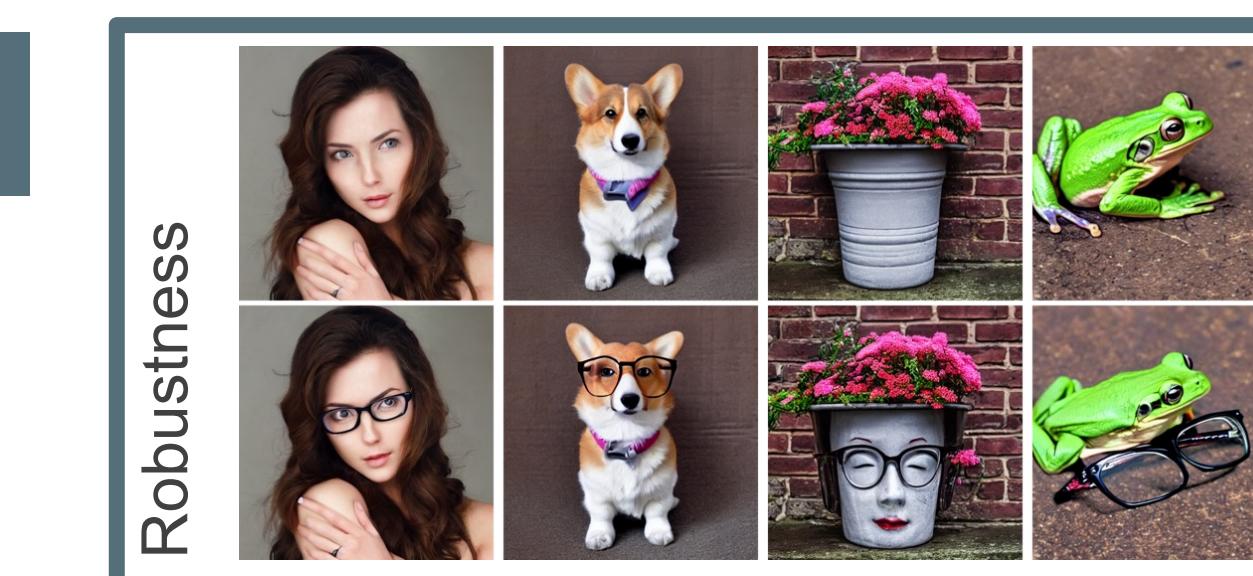
Composable
Diffusion

Properties

Isolation



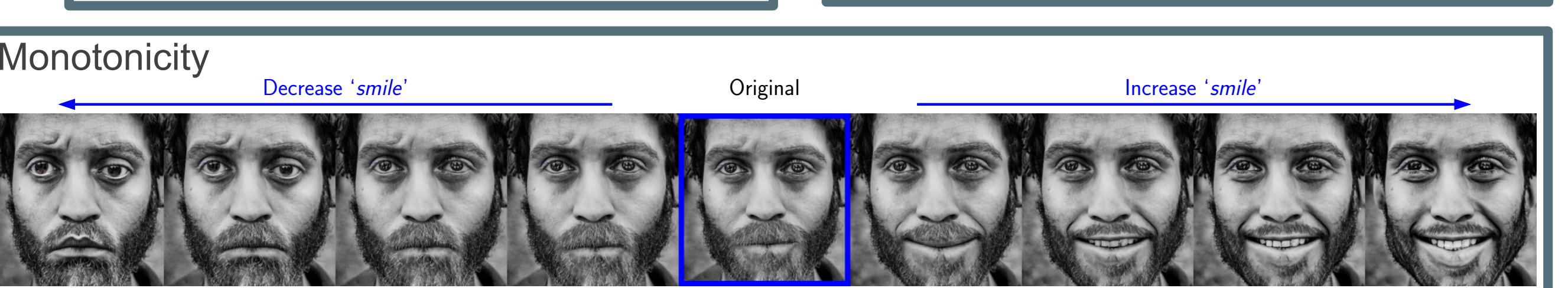
Robustness



Uniqueness

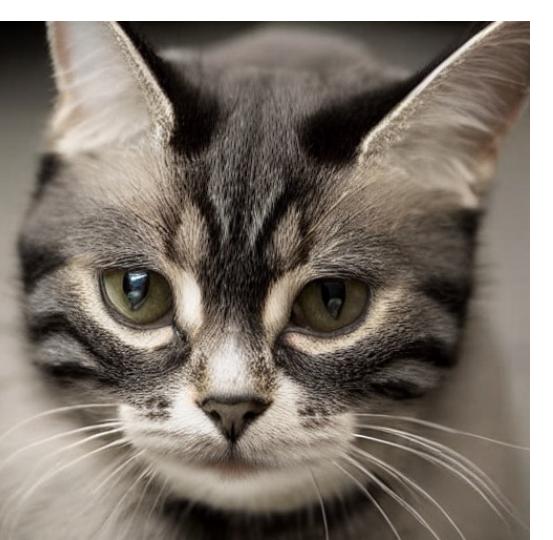


Monotonicity



Failure when changing input prompt

"A photo
of a cat"



"A photo of a cat
wearing sunglasses"



& coherent manipulation

"A photo of a cat"
+SEGA: "sunglasses"



Code

Integrated in Diffusers



Demo



Conclusion

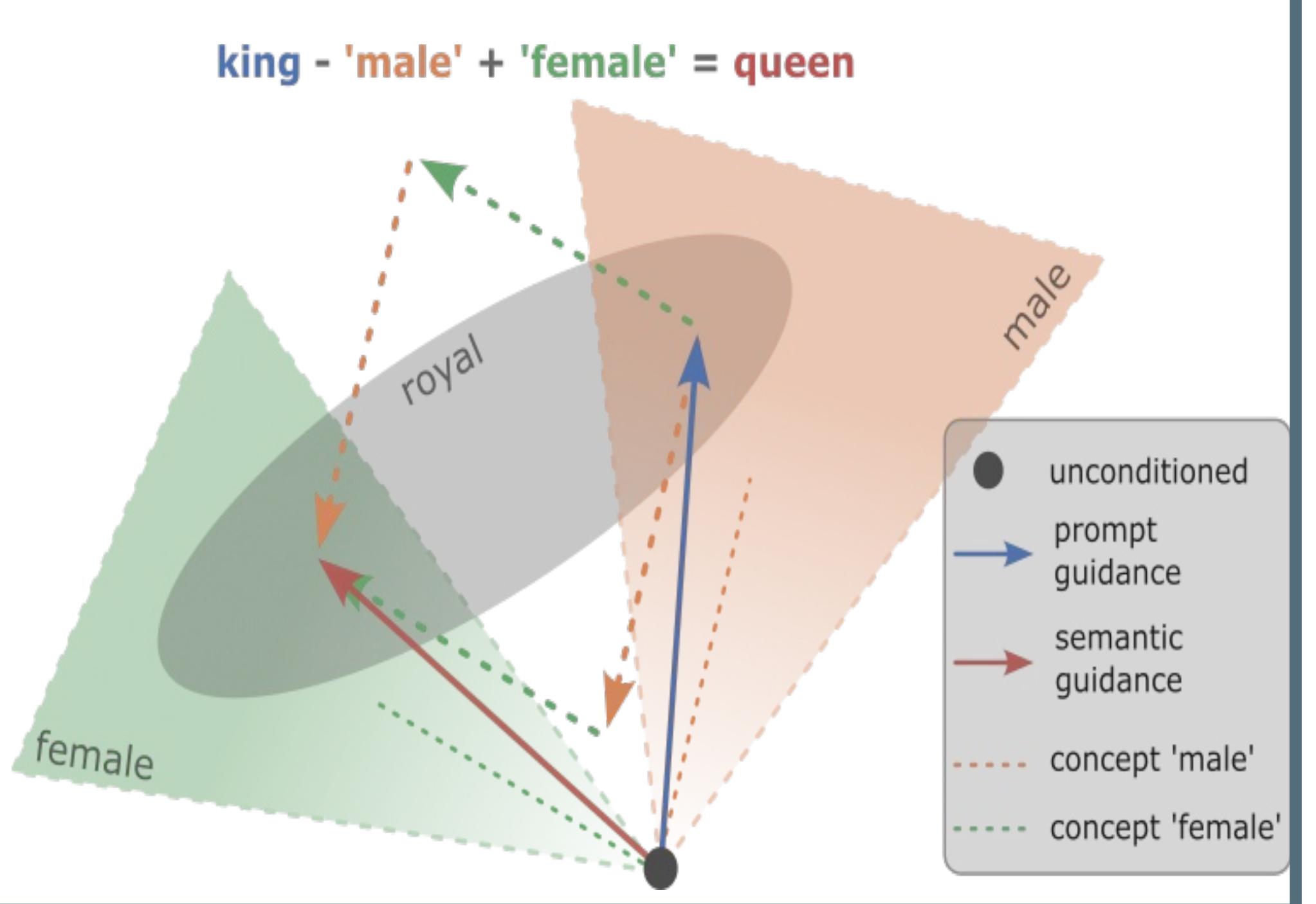
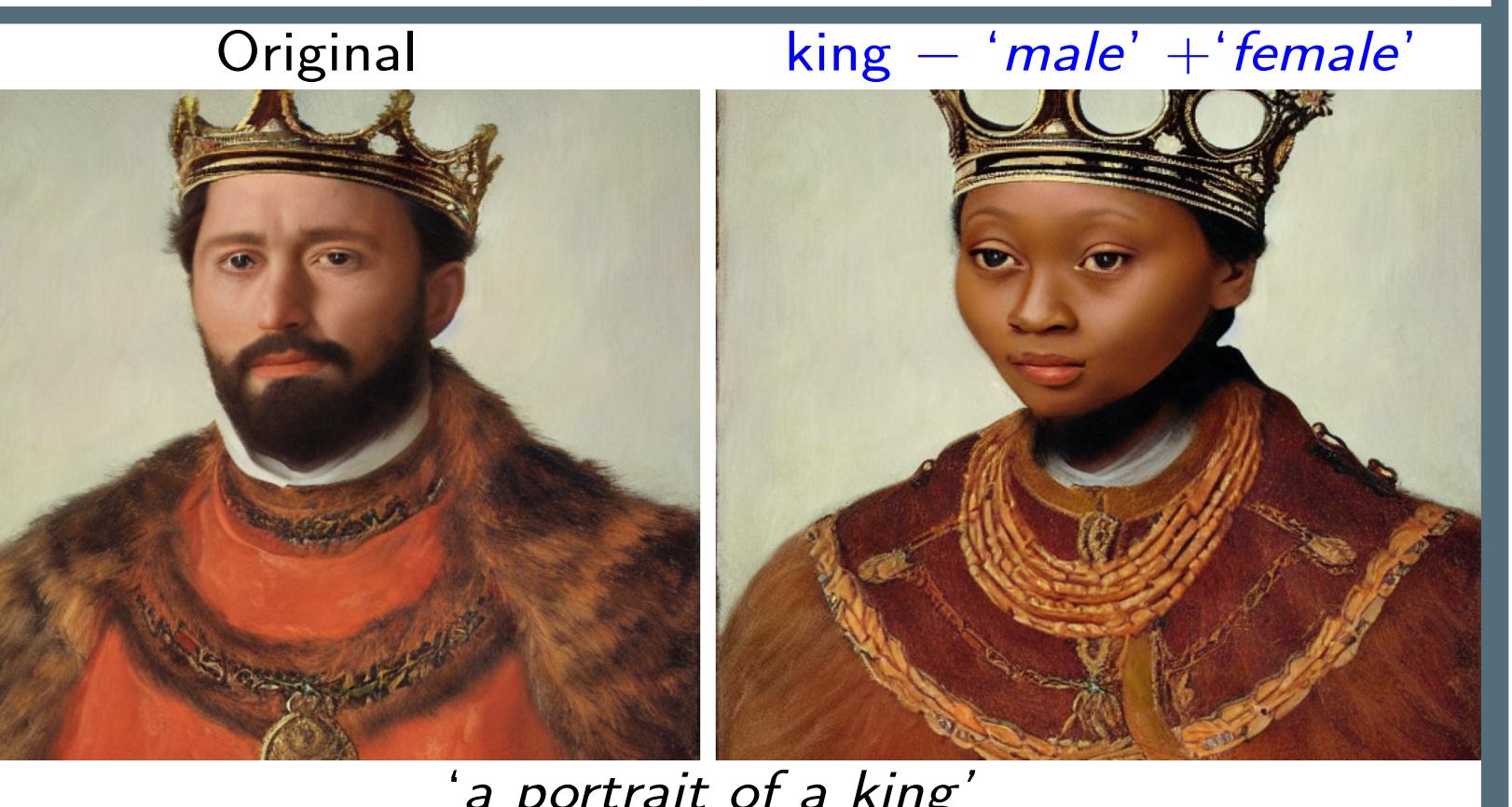
- SEGA enables **semantic control** over image generation
- at **inference**, with **no training**, and **architecture-agnostic**

Check out the Project Page

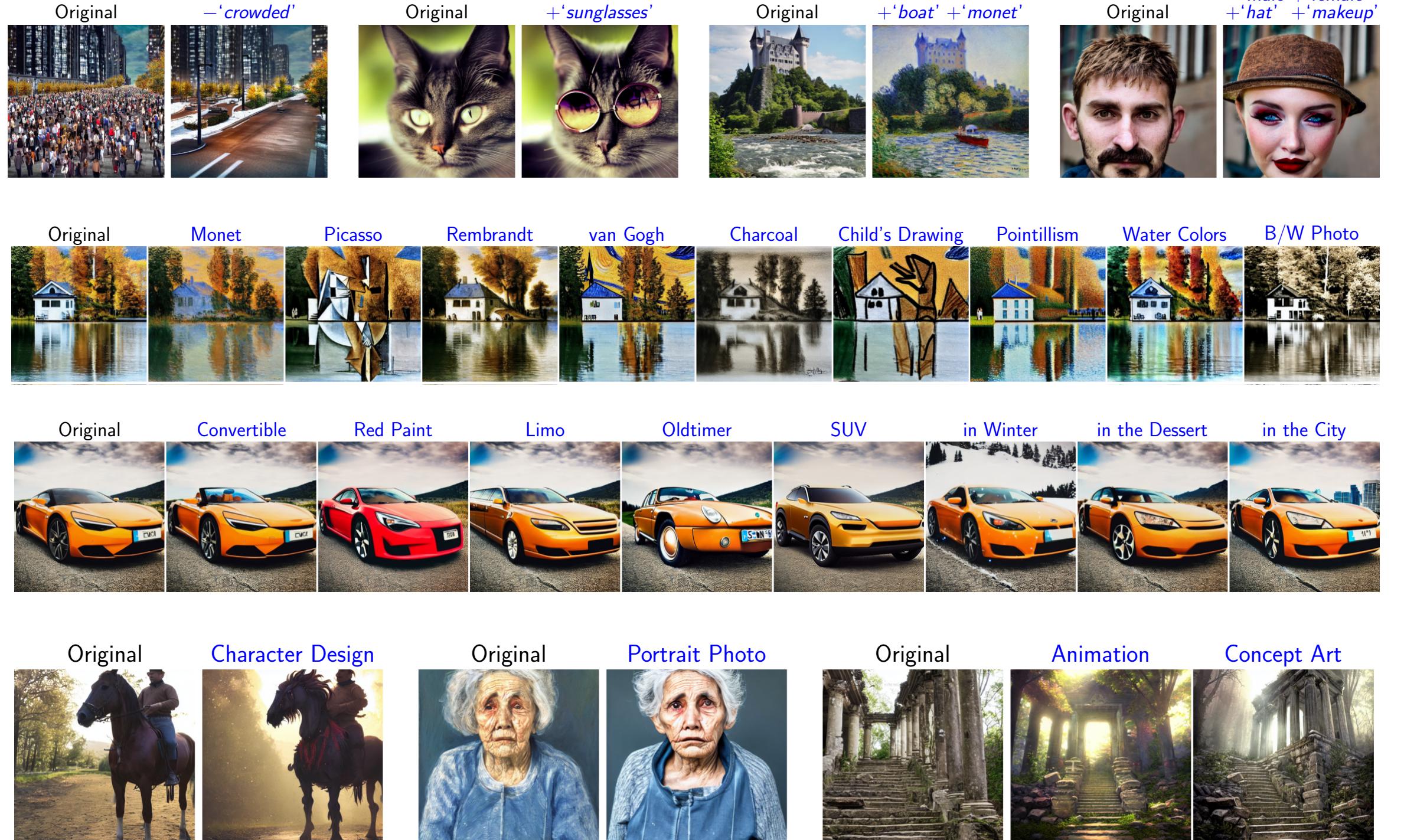
Semantic Guidance

Semantic Guidance

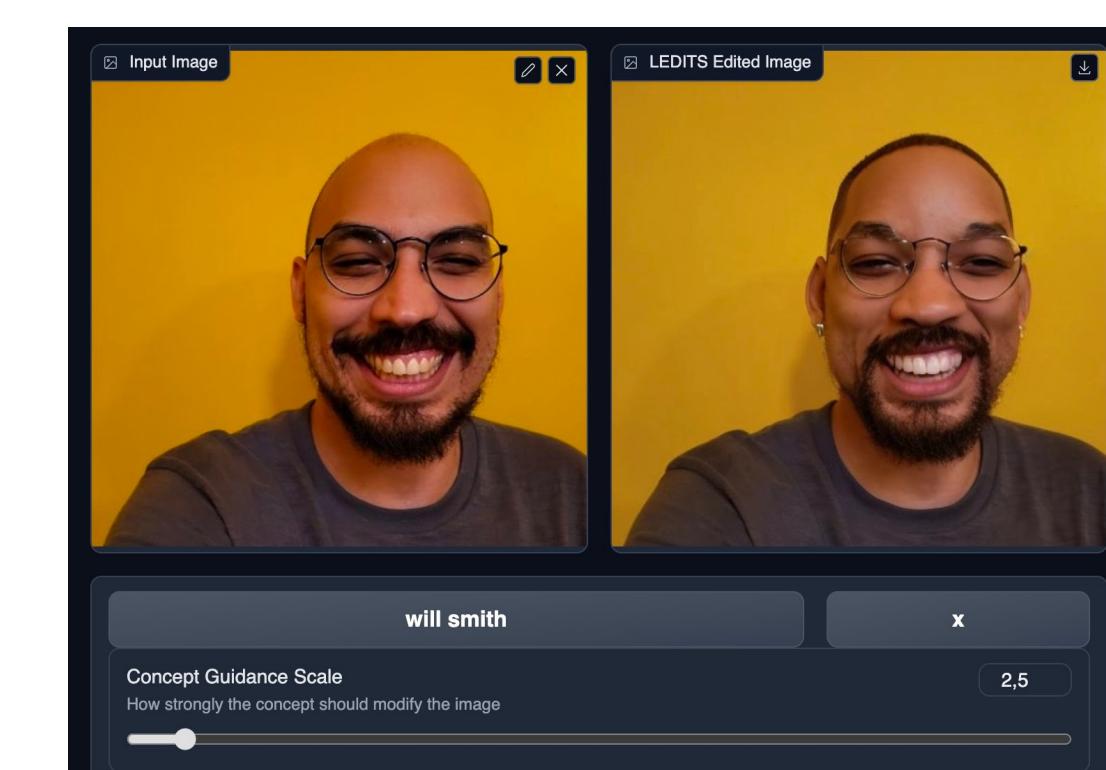
$$\begin{aligned}\bar{\epsilon}_\theta(z_t, c_p, c_e) = \\ \epsilon_\theta(z_t) + s_g(\epsilon_\theta(z_t, c_p) - \epsilon_\theta(z_t)) + \gamma(z_t, c_e) \\ \gamma(z_t, c_e) = \mu(\psi; s_e, \lambda)(\pm(\epsilon_\theta(z_t, c_e) - \epsilon_\theta(z_t)))\end{aligned}$$



Applications



Follow up work: Real Image Editing with LEdits++



Workshop on ML for Creativity and Design

- Sat 16 Dez.
- Talk: 9:20 am
- Poster: 1:30 pm

