# Versatile Diffusion: Text, Images and Variations All in One Diffusion Model

— Xingqian Xu, Zhangyang Wang, Eric Zhang, Kai Wang, Humphrey Shi

# Outline

- History of Generative Models
- Versatile Diffusion
- Network Architecture
- Disentanglement of Style and Semantic
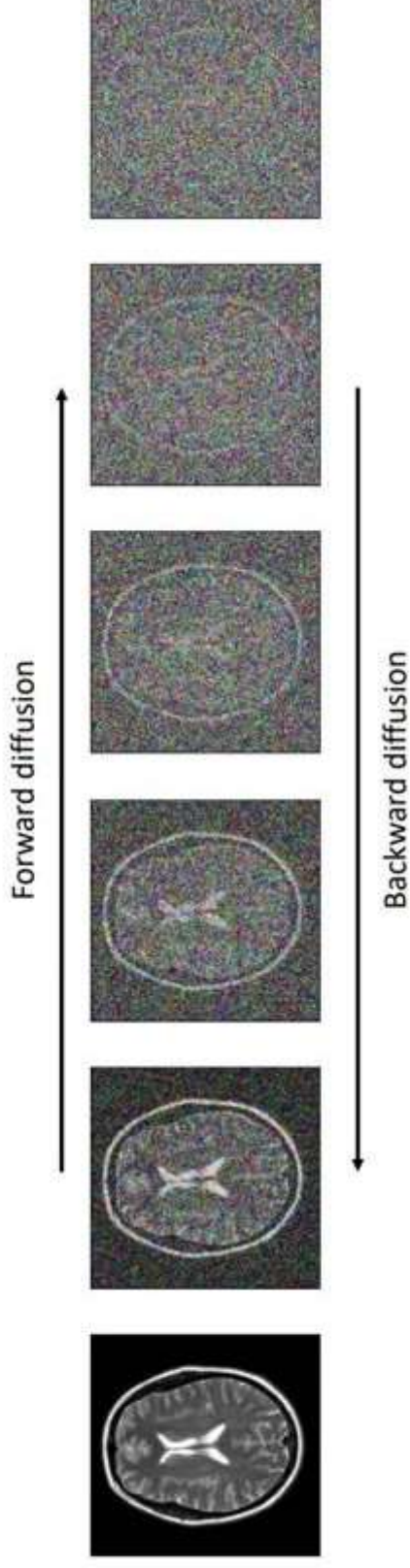- Dual Context Blender
- Questions

# Generative Adversarial Networks

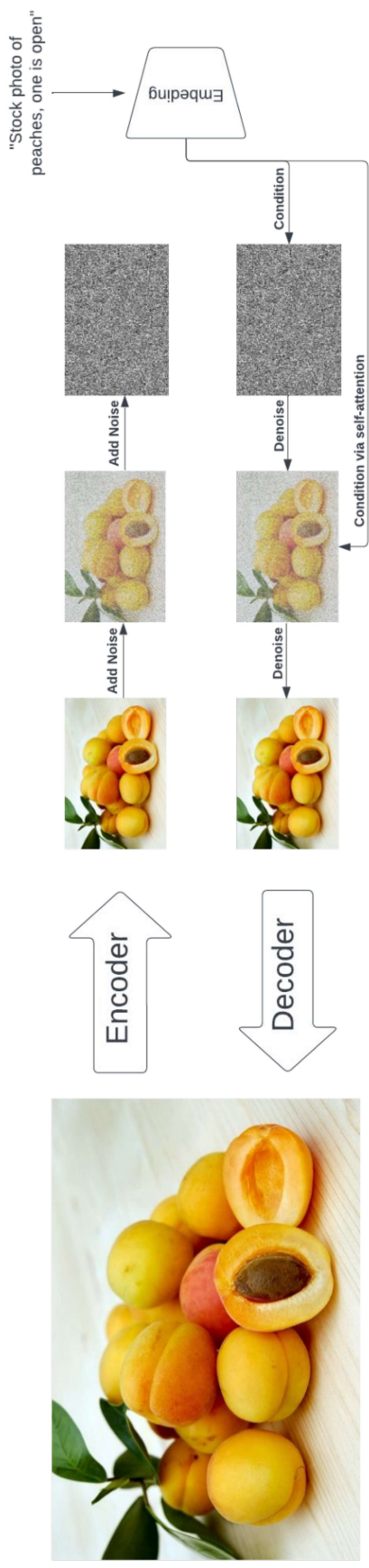GAN focus on specific domains and specific tasks i.e. faces

# New Horizons: Diffusion Model

Likelihood based models that gradually restore image contents from gaussian corruptions.
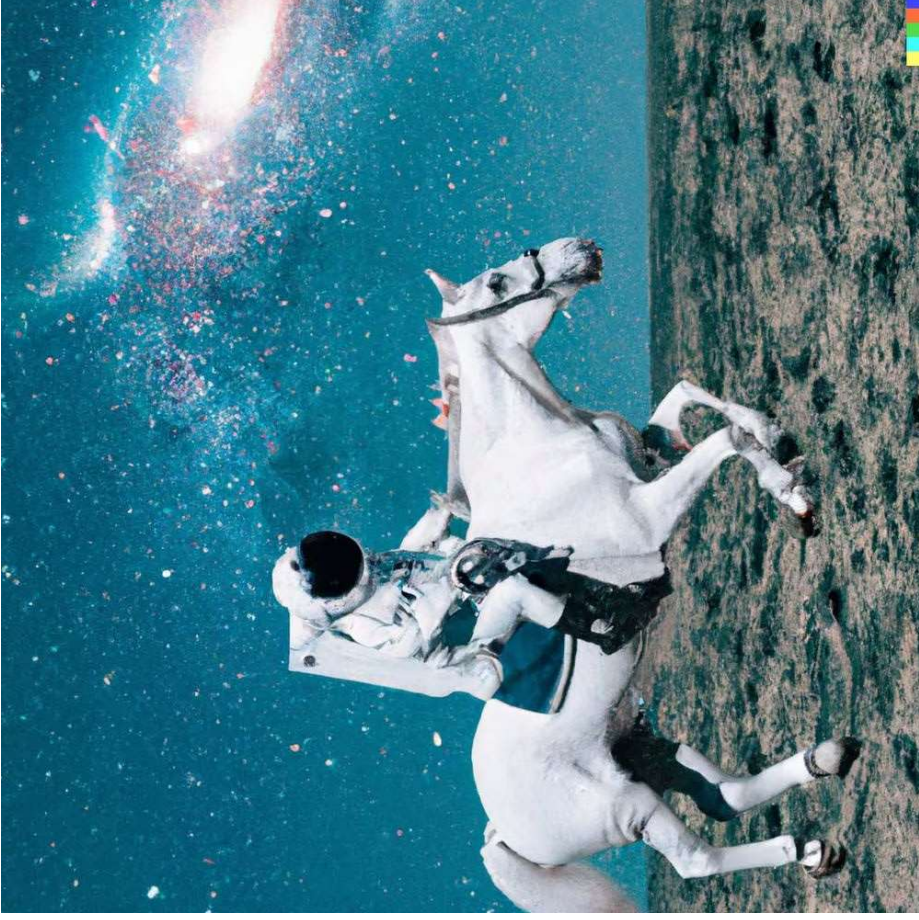


Forward diffusion

Backward diffusion

# New Horizons: Diffusion Model

It has proven to be effective in bridge modalities and tasks.

# DALL-E2 (2021)

Input: An astronaut riding a horse in photorealistic style.

## DALL-E3 (2023)



A bustling city street under the shine of a **full moon.**

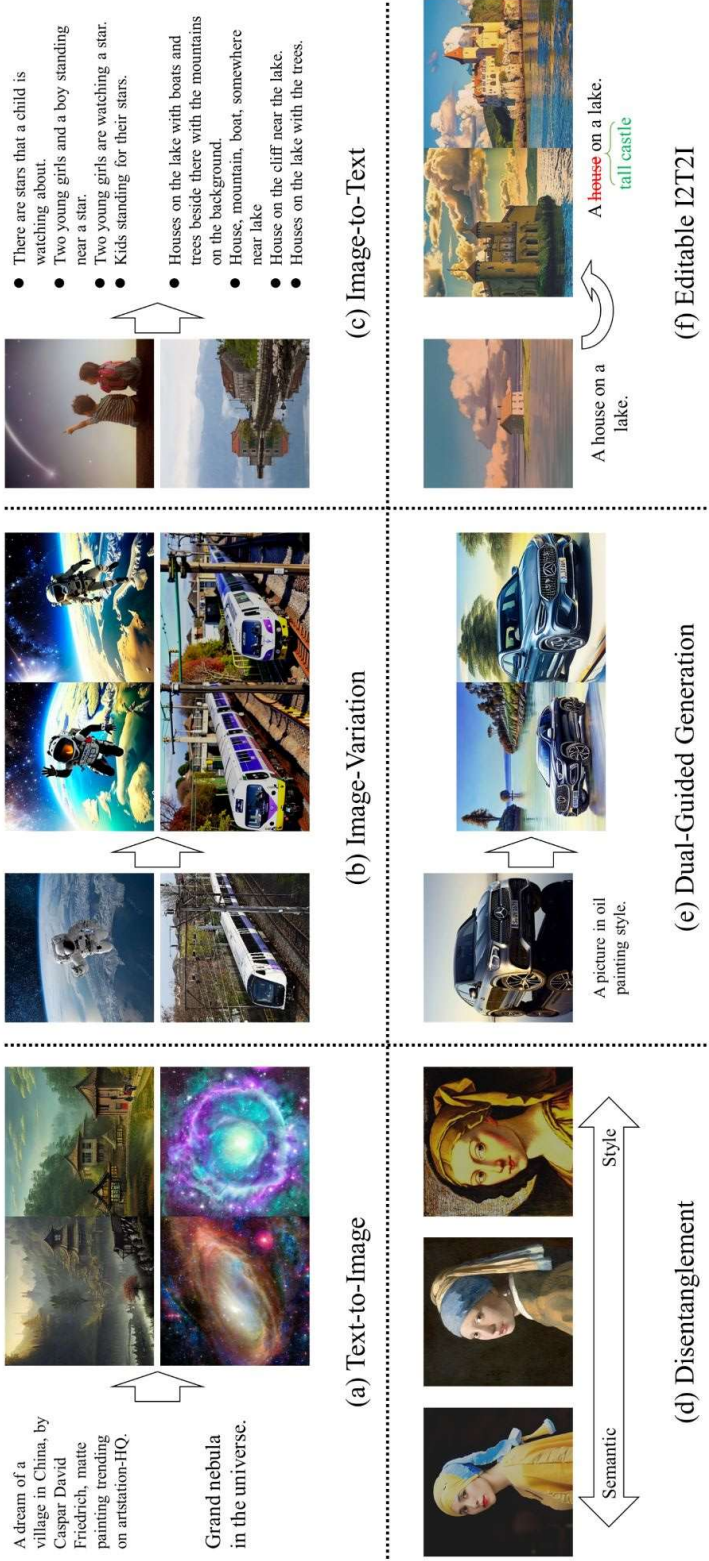The sidewalks bustling with **pedestrians enjoying the nightlife.**

The grumpy vendor, a **tall, sophisticated man,** is wearing a sharp suit, sports a **noteworthy moustache** and is animatedly conversing on his **steampunk telephone.**
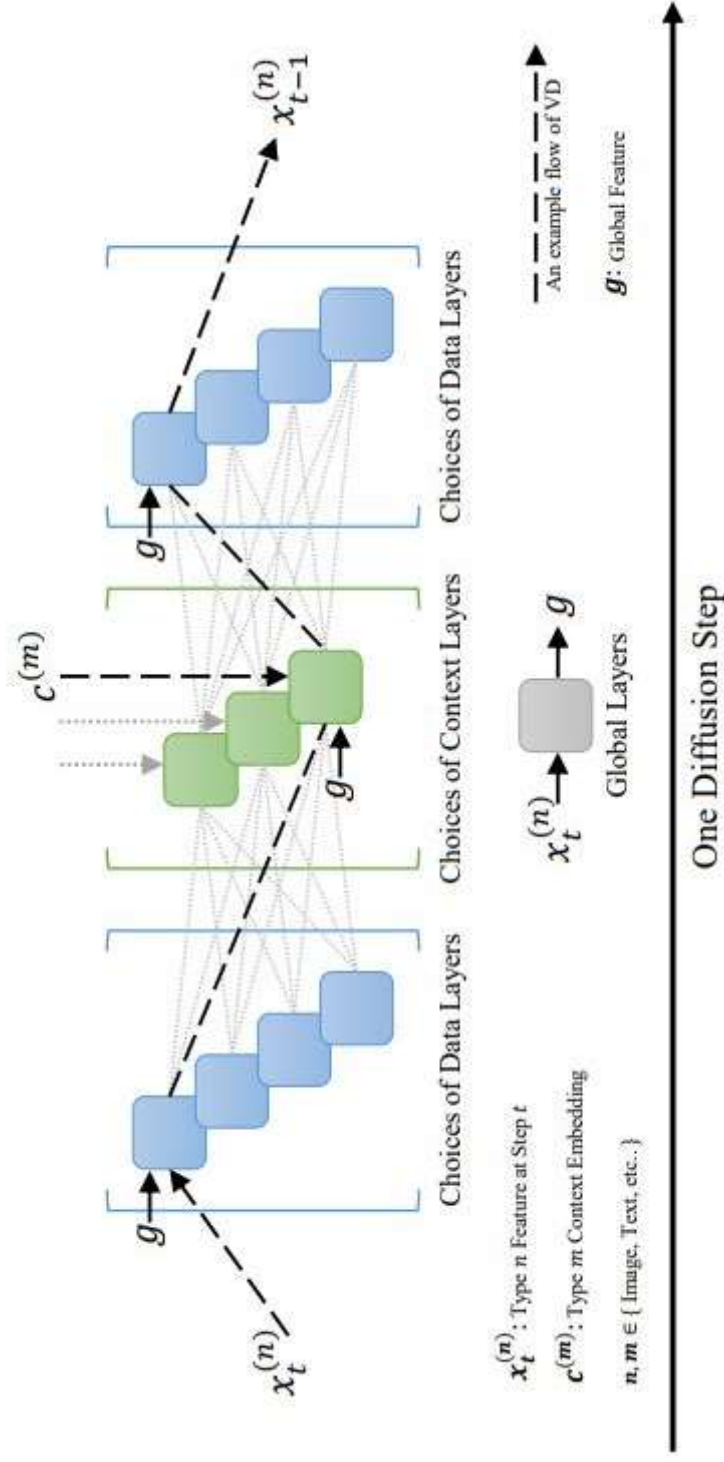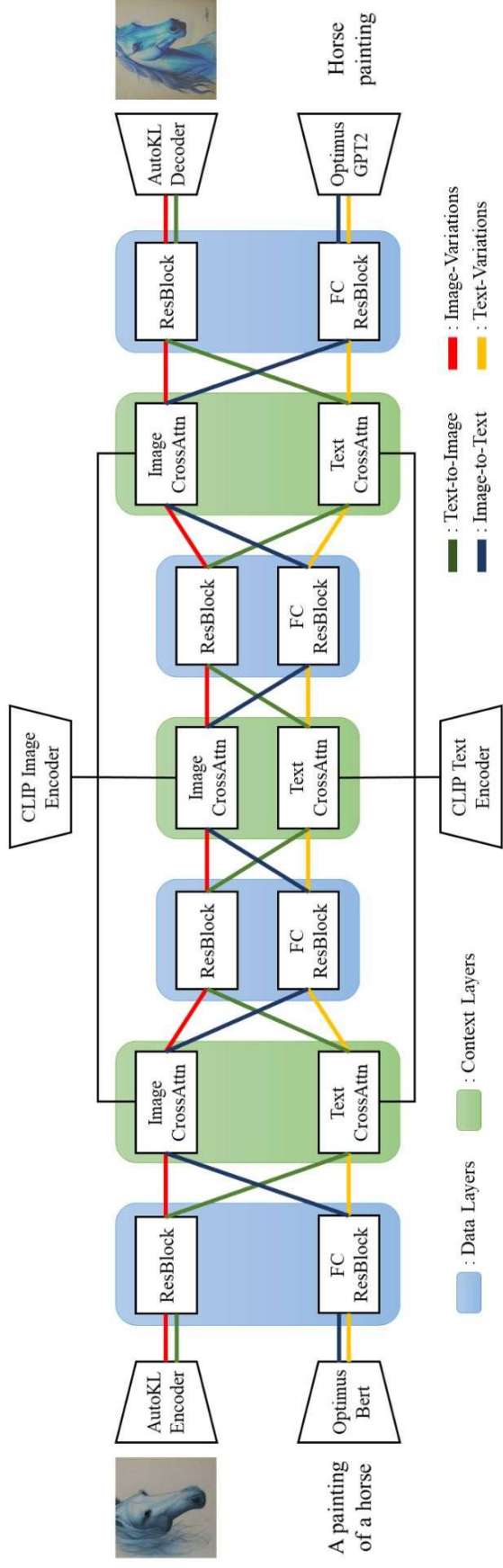
At the corner stall, a **young woman** with fiery red hair, dressed in a signature velvet cloak, is **haggling with the grumpy old vendor.**

# What's Next?

# Versatile Diffusion (2023)

A dream of a village in China, by Caspar David Friedrich, matte painting trending on artstation-HQ.

Grand nebula in the universe.

**(a) Text-to-Image**

**(b) Image-Variation**

- There are stars that a child is watching about.
- Two young girls and a boy standing near a star.
- Two young girls are watching a star.
- Kids standing for their stars.

- Houses on the lake with boats and trees beside there with the mountains on the background.
- House, mountain, boat, somewhere near lake
- House on the cliff near the lake.
- Houses on the lake with the trees.

**(c) Image-to-Text**

A picture in oil painting style.

**(e) Dual-Guided Generation**

Semantic ⟷ Style

**(d) Disentanglement**

A house on a lake.

A ~~house~~ on a lake.
tall castle

**(f) Editable I2T2I**

# One Diffusion Step



Choices of Data Layers

Choices of Context Layers

Choices of Data Layers

Global Layers

$x_t^{(n)} \rightarrow g$

One Diffusion Step

$x_t^{(n)}$ : Type $n$ Feature at Step $t$

$c^{(m)}$ : Type $m$ Context Embedding

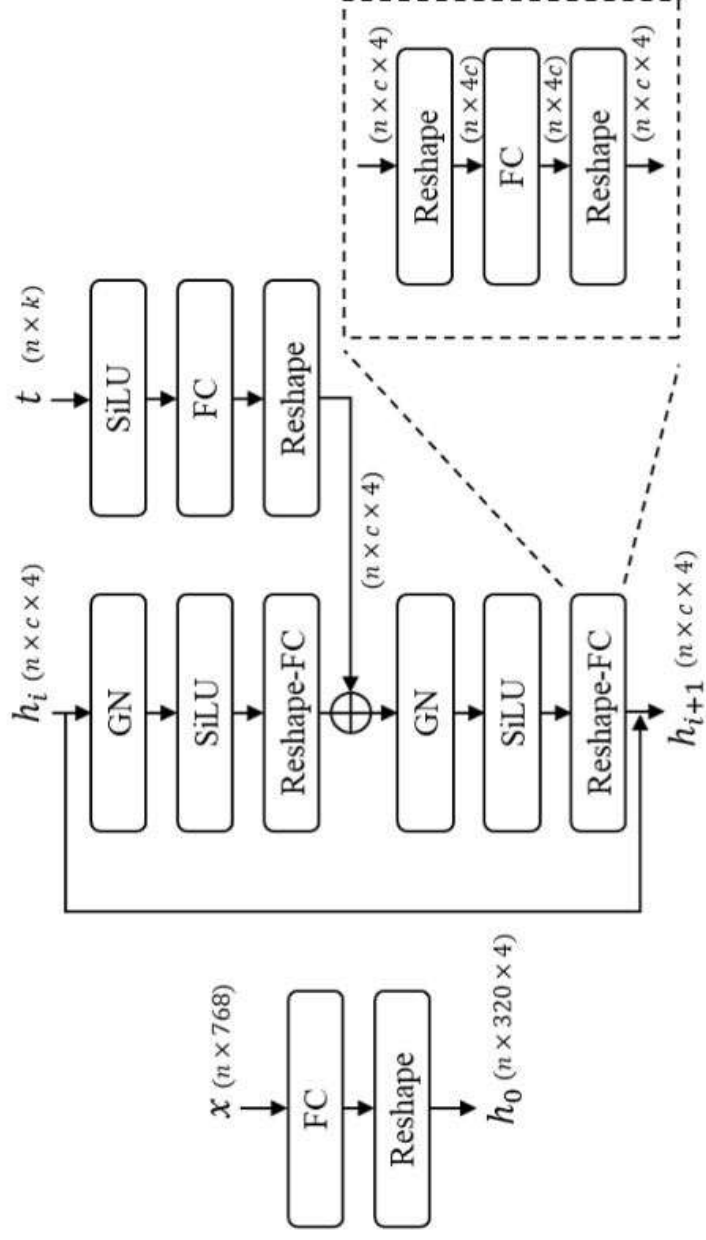$n, m \in \{$ Image, Text, etc..$\}$

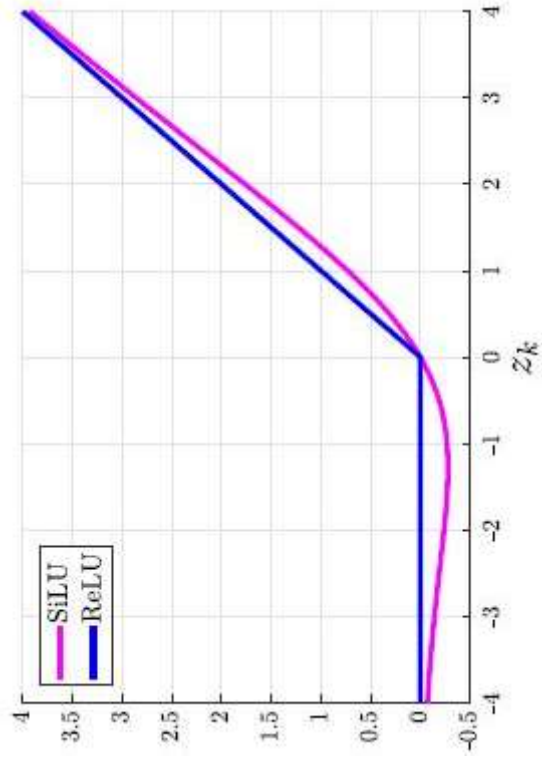- - - An example flow of VD

$g$ : Global Feature

# Network Structure
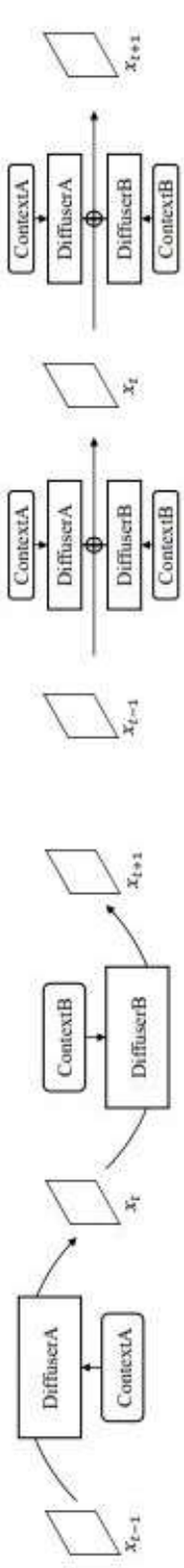
# Text Data Layers

**SiLu**

# Training

- For each of the flows compute the variational weighted losses and do regular backpropagation.

- Update model weights when the gradient in all flows are accumulated.

- VD is trained progressively in three settings:
  - Single Flow: Image Variation
  - Dual Flow: Text-to-Image and Image-variation
  - Four Flow: All the tasks together

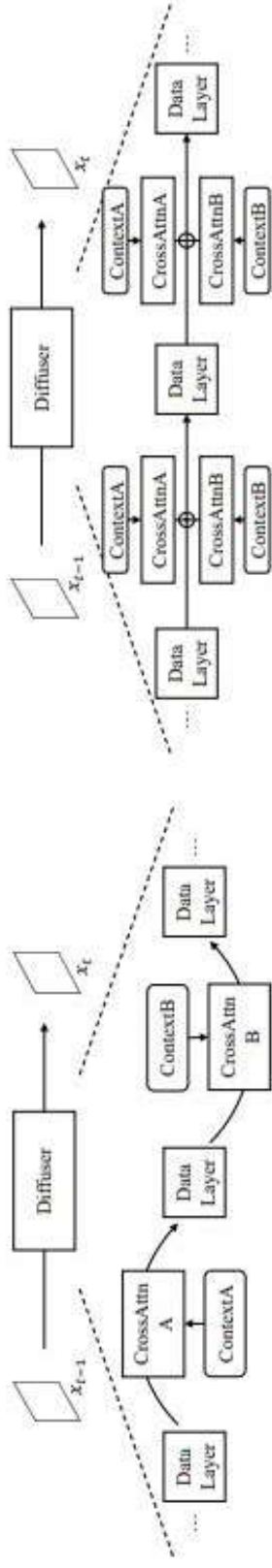# Disentanglement of Style and Semantic



Input          Semantic Focused          Variation          Style Focused

# Dual-context Blender



(a) Model-level Mixing A
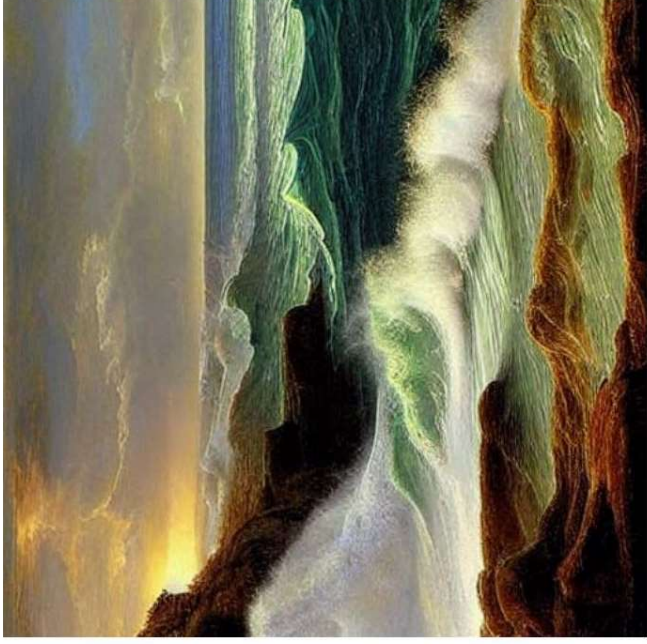
(b) Model-level Mixing B

(c) Layer-level Mixing

(d) Attention-level Mixing

# Results: Text to Image



"A wonderful evening in New York City with a great view of Brooklyn Bridge and a magnificent city view of Manhattan, HD 8K"



"A beautiful painting of waves crashing on a cliff by Thomas Cole"

# Results: Image Variation



Input



Variation #1



Variation #2

# Results: Dual-Context Blender



Input



"100 mph"



"Traveling among the stars"

# Limitations

- **Limited Latent Space**: Optimus VAE's latent vector are 768 single dimension generated using Bert which might be inadequate for long text sentences. It is weak in understanding word locations and orders

- **Imperfect Text Data**: There is domain shift in Optimus VAE's training data compare to VD's training data making it difficult to reconstruct certain images.

# Question 1

How does the authors disentangle and semantic using the VD?

# Question 2

Give three example of basic tasks VD can achieve?