

DEADiff: An Efficient Stylization Diffusion Model with Disentangled Representations

Tianhao Qi^{1*}, Shancheng Fang¹, Yanze Wu^{2†}, Hongtao Xie^{1✉}, Jiawei Liu²,
Lang Chen², Qian He², Yongdong Zhang¹

¹University of Science and Technology of China, ²ByteDance,

*Works done during the internship at ByteDance, †Project lead, ✉Corresponding author,

PRESENTER: GUO TANG

2024/5/19

■ Outline

1 / **Background**

2 / Author

3 / Method

4 / Experiments

■ Background

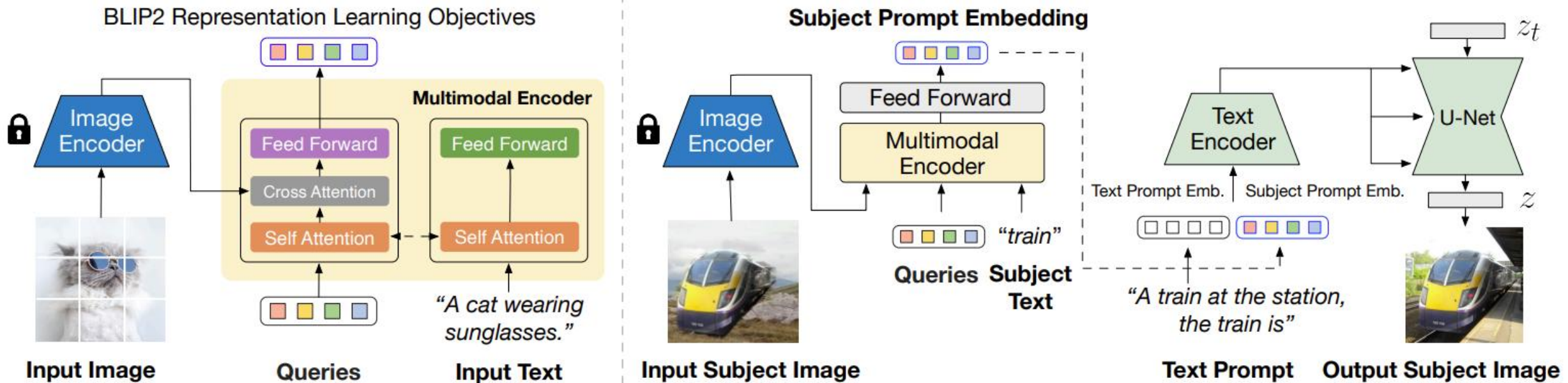
BLIP-Diffusion: Pre-trained Subject Representation for Controllable Text-to-Image Generation and Editing

Dongxu Li[†], Junnan Li[†], Steven C.H. Hoi[†]
Salesforce AI Research

[†]Corresponding authors: {li.d, junnan.li, shoi}@salesforce.com
<https://github.com/salesforce/LAVIS/tree/main/projects/blip-diffusion>

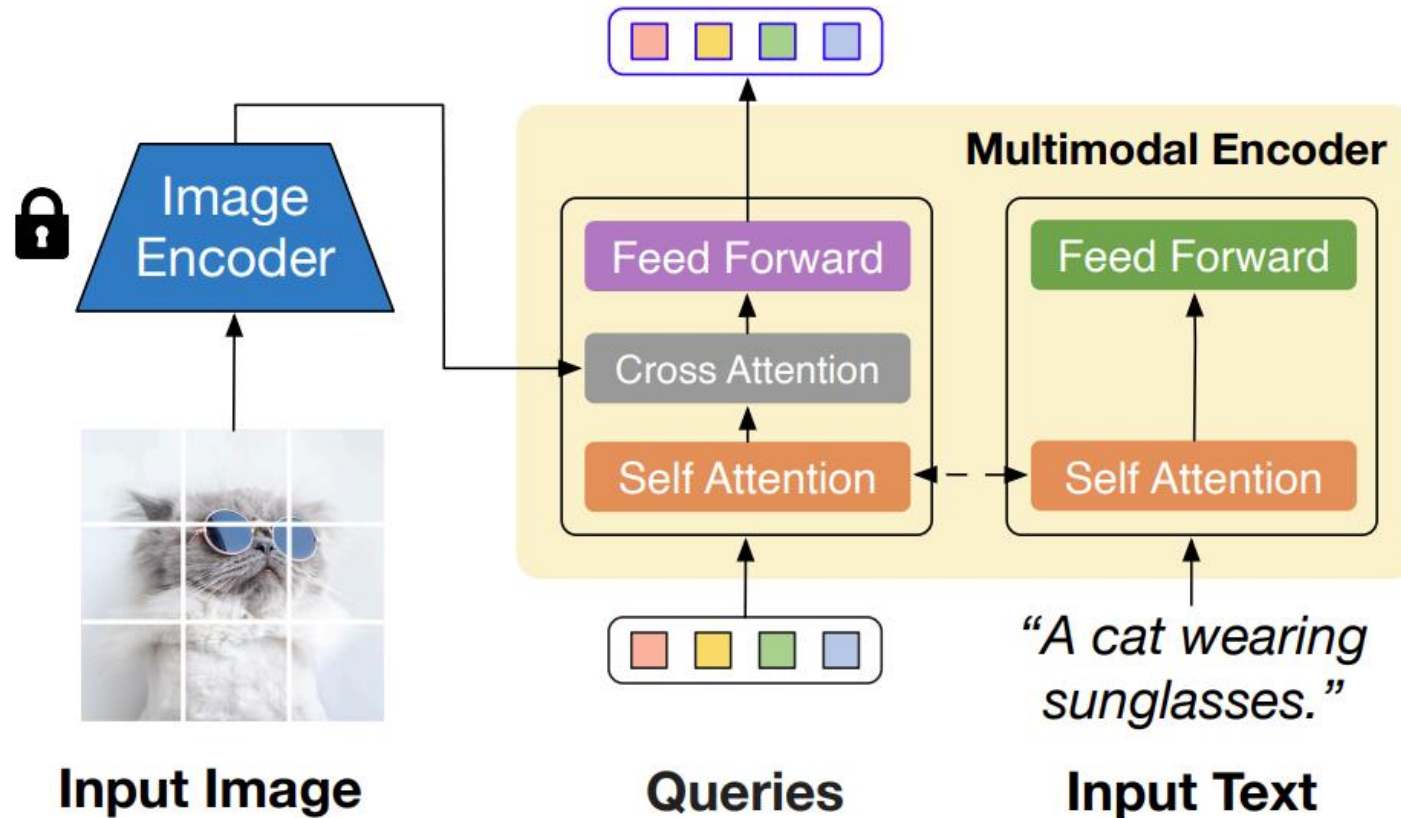
Background

Two stage pre-training



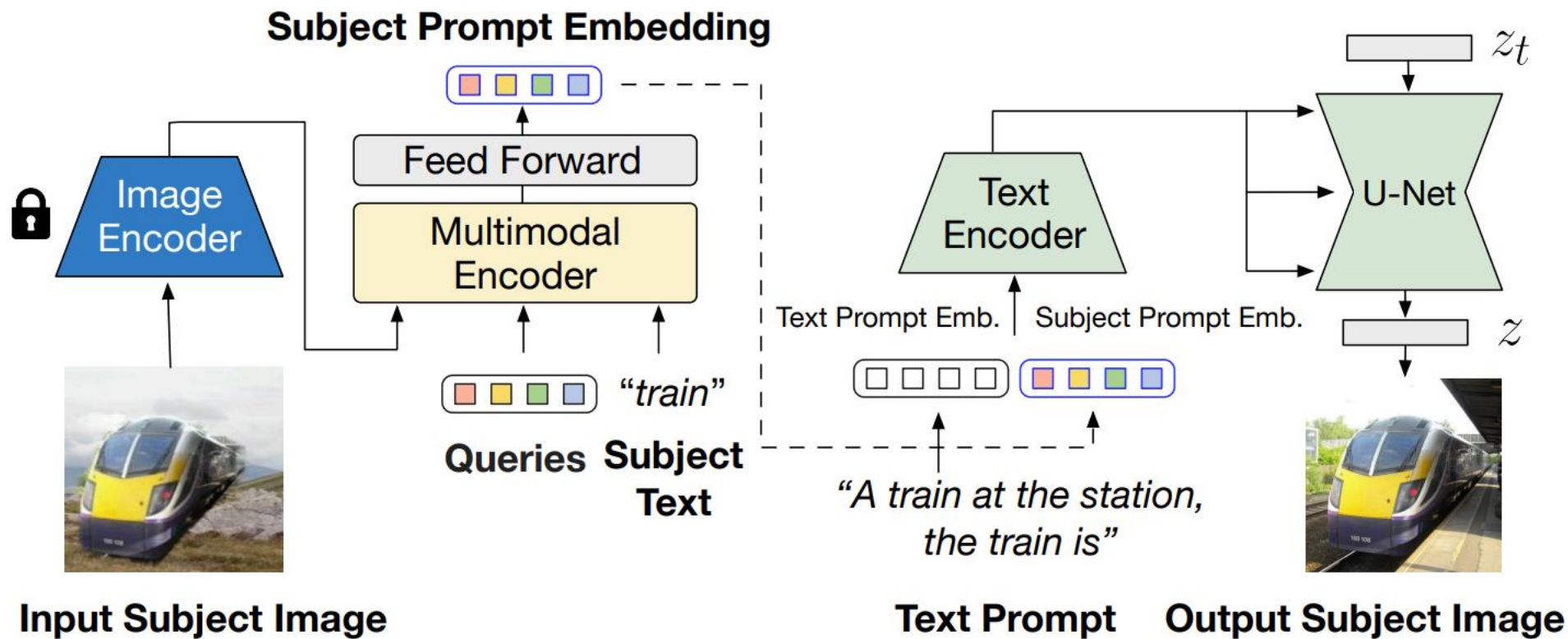
■ Background

BLIP2 Representation Learning Objectives



Background

BLIP2 Representation Learning Objectives



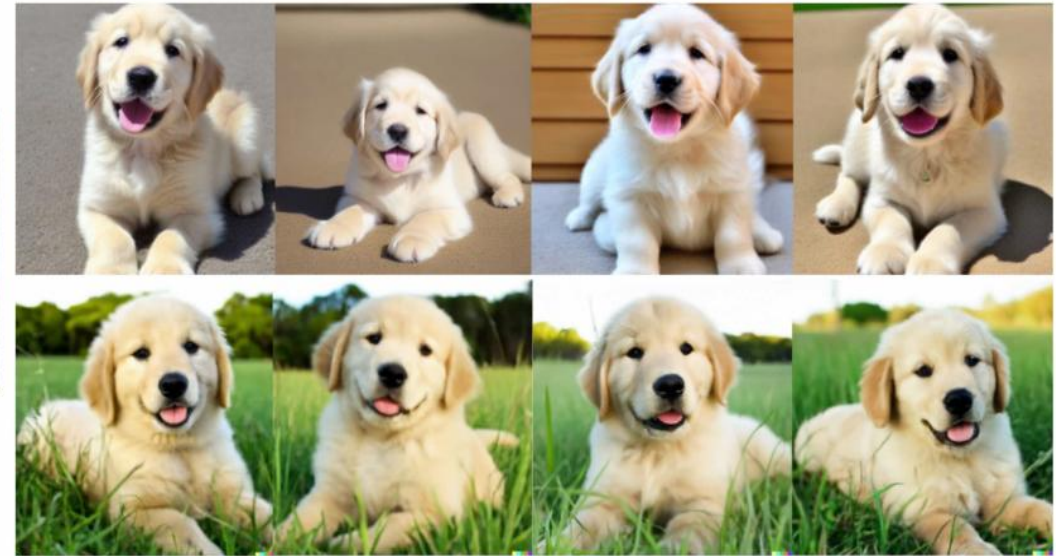
[text prompt], the [subject text] is [subject prompt]"

■ Background

Generating training image pair

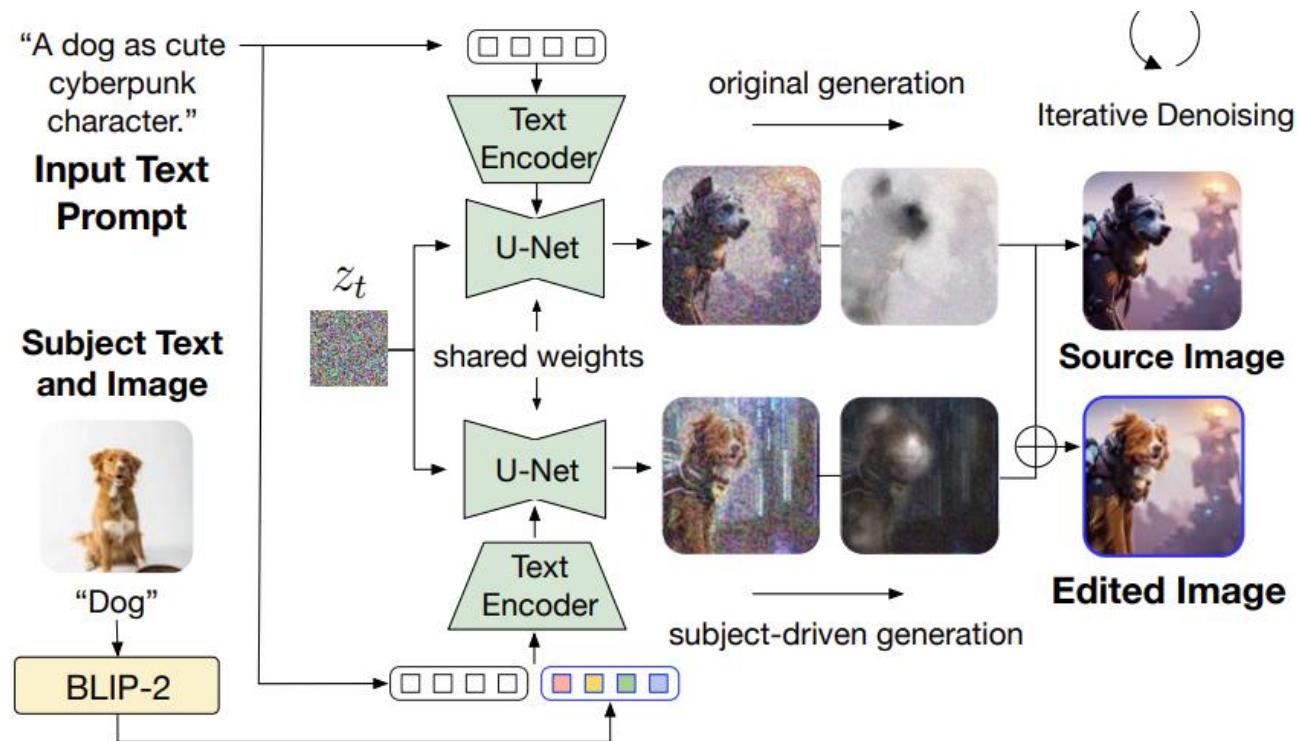
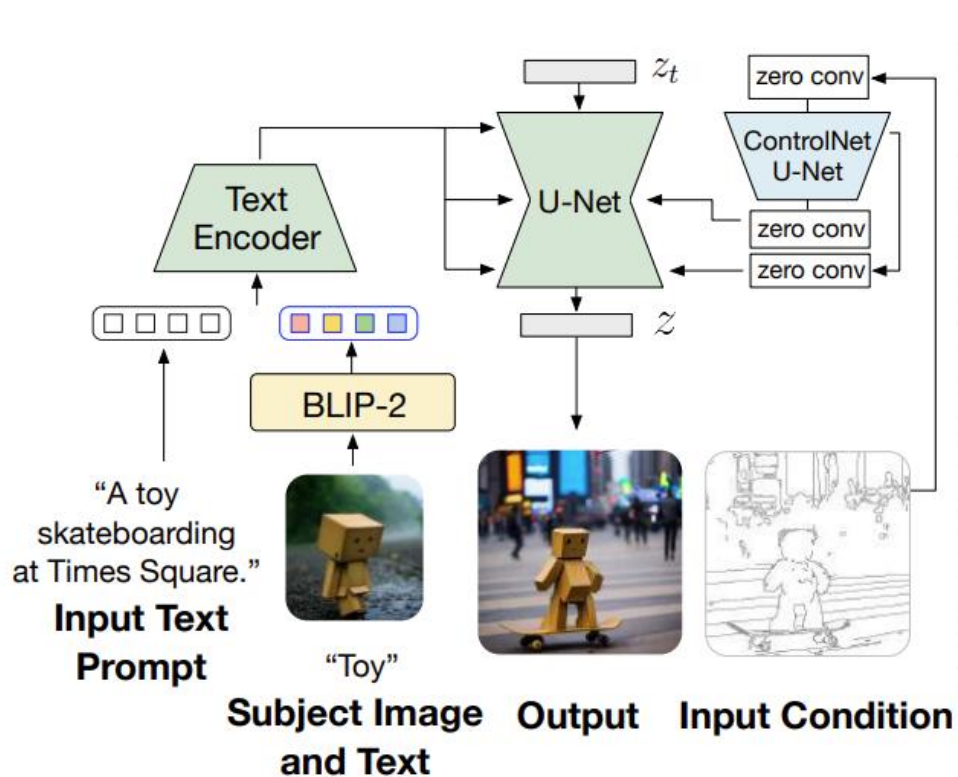


Input image



Background

- Subject-specific Fine-tuning and Inference
- Structure-controlled Generation with ControlNet
- Subject-driven Editing with Attention Control



■ Background

zero-sort subject-driven generation



Input Image

cat



on a skateboard



as plushie



Input Image

car



painting by Van Gogh



colored green



Input Image

toy



in a forest



with sunset



wearing top hat



on a red rug



at Palace of Versailles



seen from back



made of lego



at grand canyon



on the beach

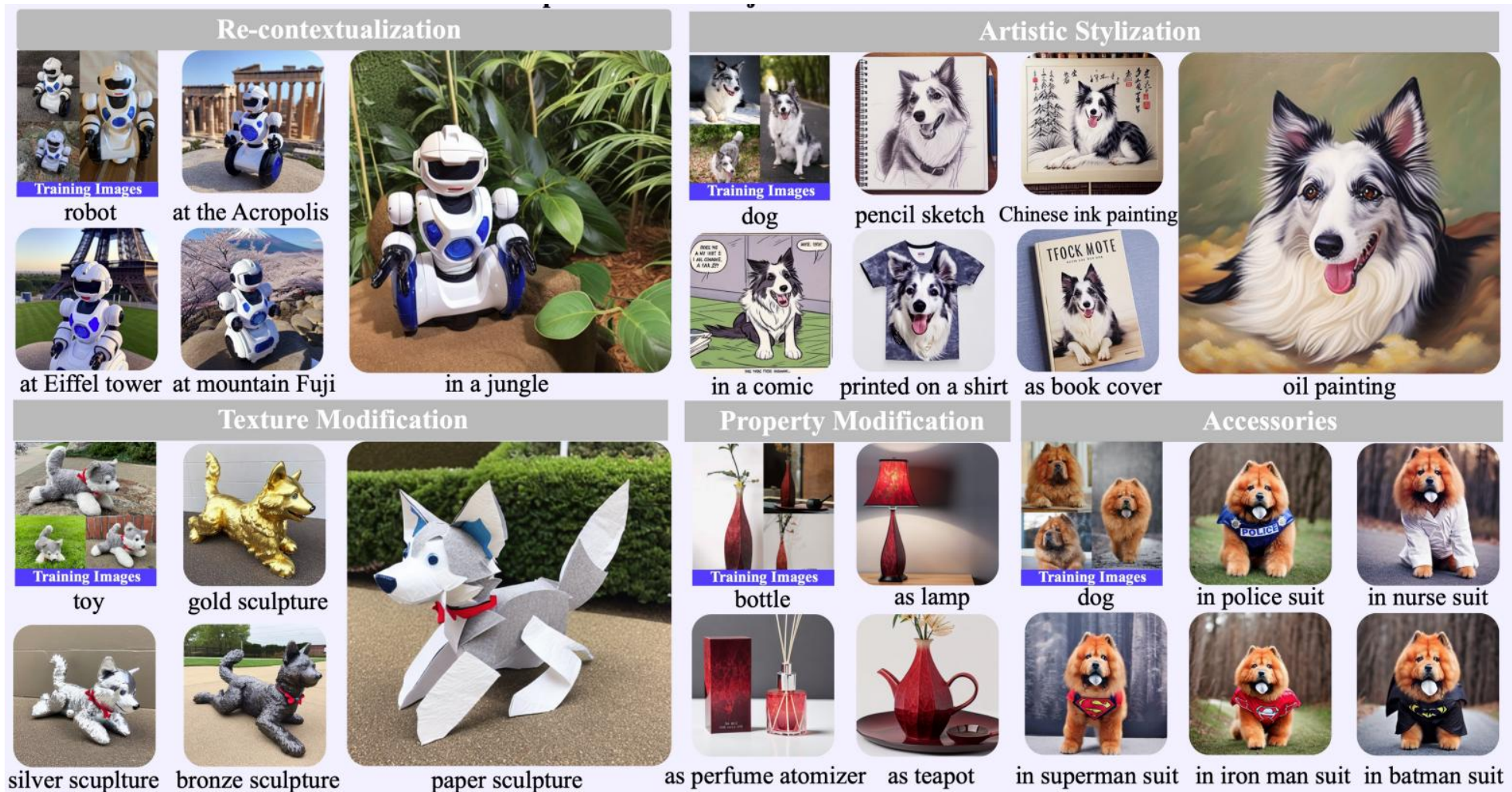


decorated on cup



wearing glasses

Background few-step fine-tuned subject-driven generation



■ Background

structure controlled subject-driven generation



Training Image

teapot



depth map



a sofa



Training Image

toy



canny edges



skateboarding at times square



Training Image

teapot



HED edges



pouring milk in to a cup

■ Background

subject-driven image editing with attention control



■ Background

- context-appearance entanglement
- failing to address the text prompt
- wrong spatial composition



Subject images



A backpack on top of a **purple rug** in a forest.



Subject images



A **cube-shaped** bear plushie.



Subject images



A sneaker **on top of** a mirror.

■ Outline

1 / Background

2 / Author

3 / Method

4 / Experiments

■ Author

DEADiff: An Efficient Stylization Diffusion Model with Disentangled Representations

Tianhao Qi^{1*}, Shancheng Fang¹, Yanze Wu^{2†}, Hongtao Xie^{1✉}, Jiawei Liu²,
Lang Chen², Qian He², Yongdong Zhang¹

¹University of Science and Technology of China, ²ByteDance,

*Works done during the internship at ByteDance, †Project lead, ✉Corresponding author,

■ Author



Tianhao Qi

PhD, University of Science and Technology of China

在 mail.ustc.edu.cn 的电子邮件经过验证

cross-modal generation object detection

First author: PhD student in University of Science and Technology of China, major in cross-modal generation, object detection

■ Author



Zhang Yongdong

University of Science and Technology of China
在 ustc.edu.cn 的电子邮件经过验证

1999-2002 Tianjin University, PhD, Signal and Information Processing

2002-2017 Researcher, Institute of Computing Technology, Chinese Academy of Sciences

Professor at the University of Science and Technology of China from 2017 to present

Research directions include multimedia content analysis, cybersecurity, and computational imaging.

Method

Results preview



Reference

(a) T2I-Adapter

(b) DEADiff

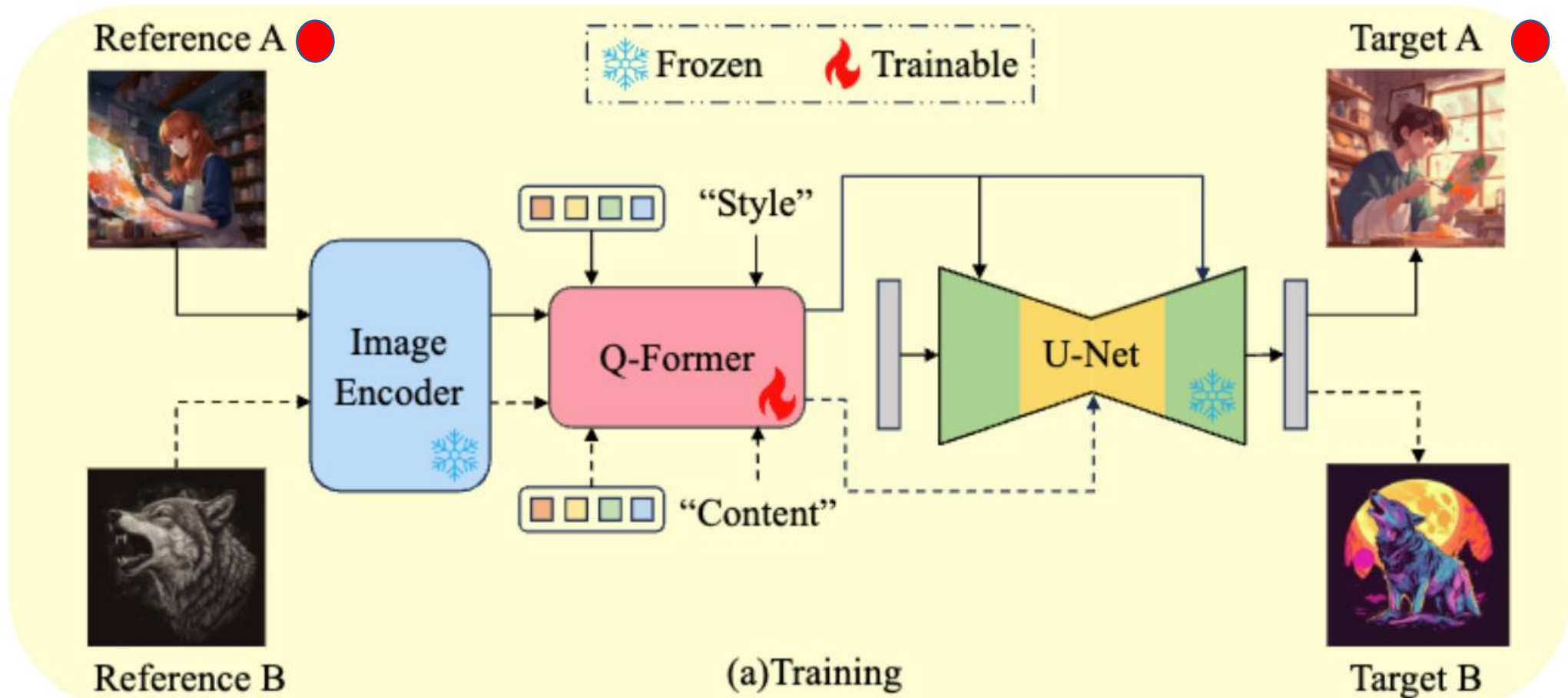
■ Method

■ Contributions

- Disentangle style and semantic representation of the reference image
- Injecting image style/semantic representation to different crossattention layers
- Established paired datasets

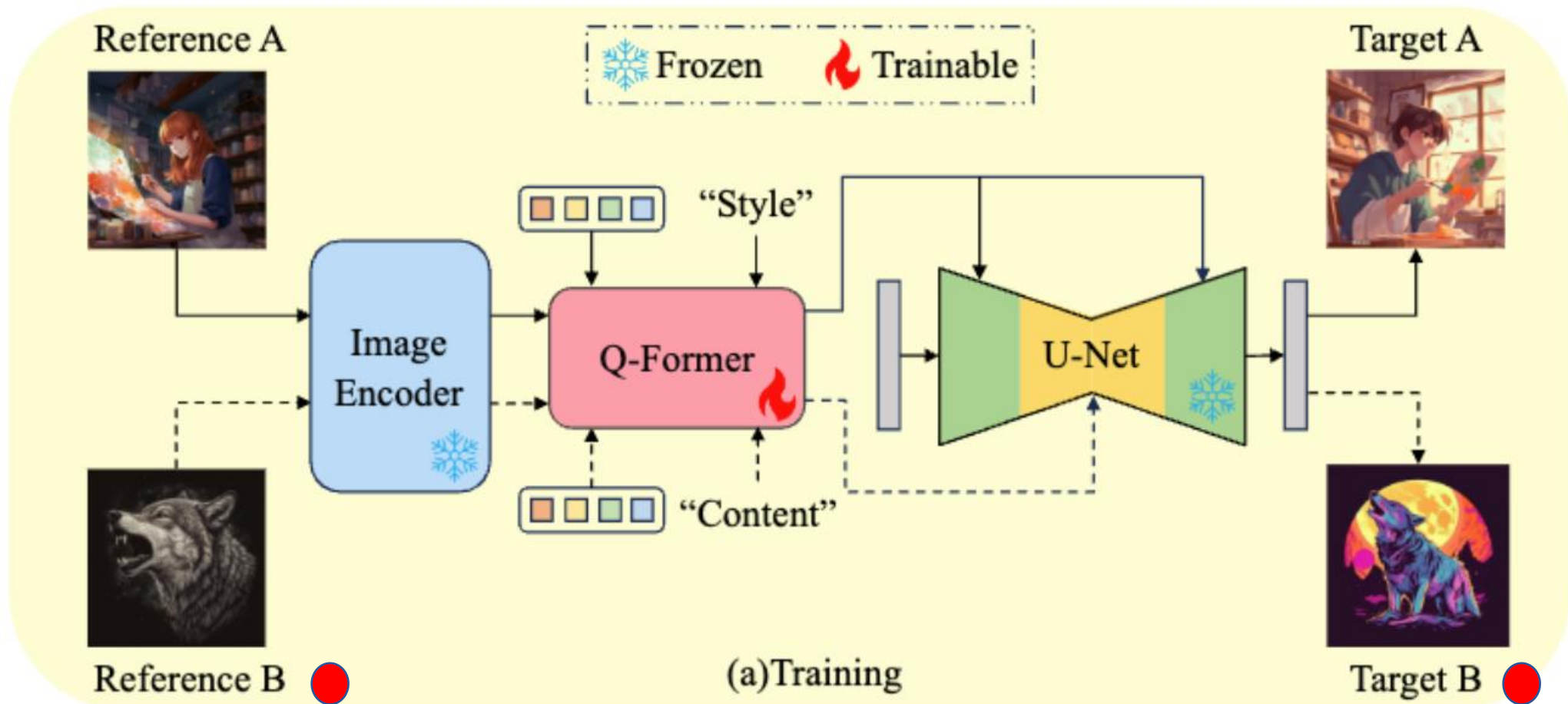
Method

- Disentangle style and semantic representation of the reference image



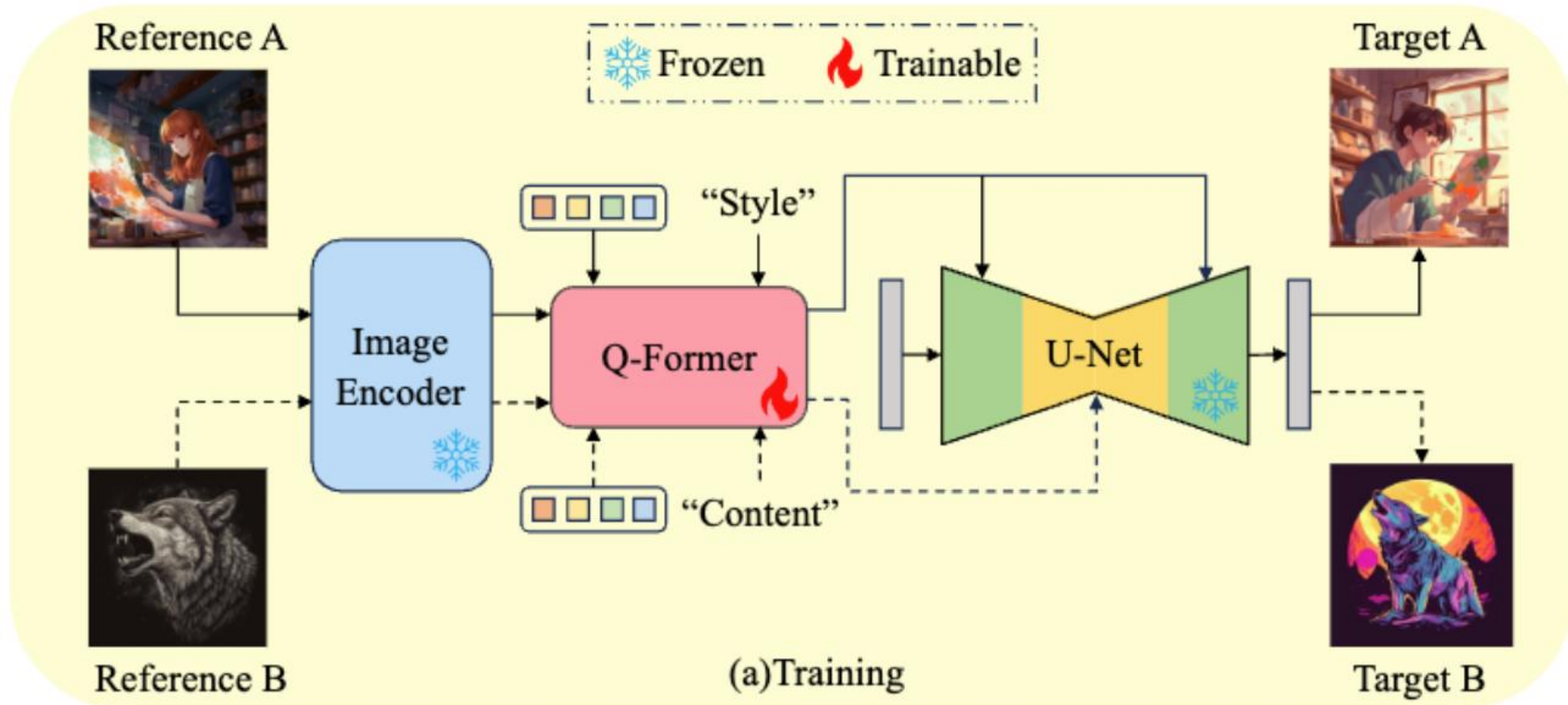
Method

- Injecting image style/semantic representation to different



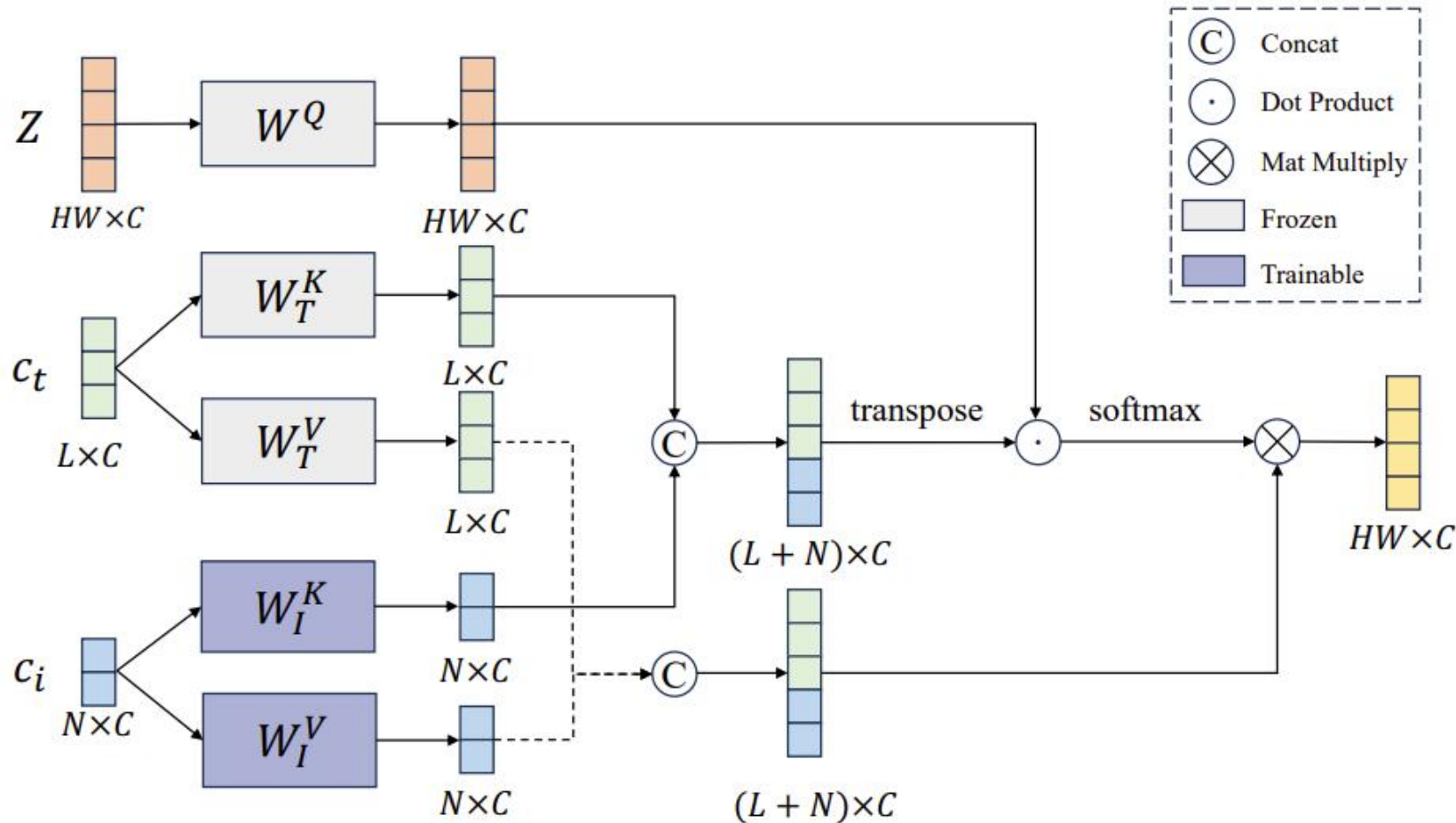
Method

- Injecting image style/semantic representation to different crossattention layers



Method

Injecting image style/semantic representation to different crossattention layers



$$Q = ZW^Q,$$

$$K = \text{Concat}(c_t W_T^K, c_i W_I^K),$$

$$V = \text{Concat}(c_t W_T^V, c_i W_I^V),$$

$$Z^{new} = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V.$$

■ Method

- Establishing paired datasets
 - Text prompt combination
 - Image generation and collection
 - Paired images selection

■ Outline

1 / Background

2 / Author

3 / Method

4 / **Experiments**

■ Experiments

- Evaluation
 - Style Similarity (SS)
 - Text Alignment capability (TA)
 - Image Quality (IQ)
 - Subjective Preference (SP)



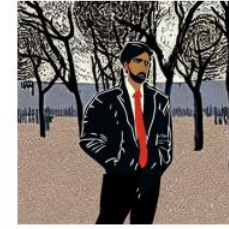
"A monkey playing with a banana"



"A palm tree"



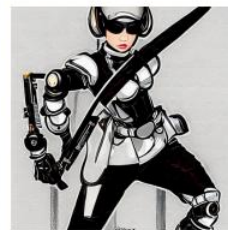
"A man wearing a black leather jacket and a red tie"



"A curly-haired boy"



"A robot"



Reference

Prompt

DEADiff

IP-Adapter
(weight tuning)

StyleAdapter

T2I-Adapter
(weight tuning)

CAST

StyleTr²

InST

■ Experiments

Quantitative Comparisons

Method	SS↑	IQ↑	TA↑	SP↑
InST [37]	0.215	5.148	0.237	6.3
CAST [36]	0.224	4.922	<u>0.282</u>	8.7
StyTr ² [3]	0.214	5.037	<u>0.282</u>	<u>13.1</u>
T2I-Adapter [17]	0.241	5.500	0.224	2.7
IP-Adapter [34]	0.274	<u>5.598</u>	0.155	-
DEADiff	<u>0.229</u>	5.840	0.284	69.0

■ Experiments

Ablations

Method	Style Similarity \uparrow	Text Alignment \uparrow
Baseline	0.274	0.148
+ DCM	0.259	0.224
+ STRE	0.222	0.286
+ SERE	0.221	0.287
DEADiff	0.224	0.289

"A dog in a bucket."



"A cactus wearing a hat."



Reference

Baseline

+DCM

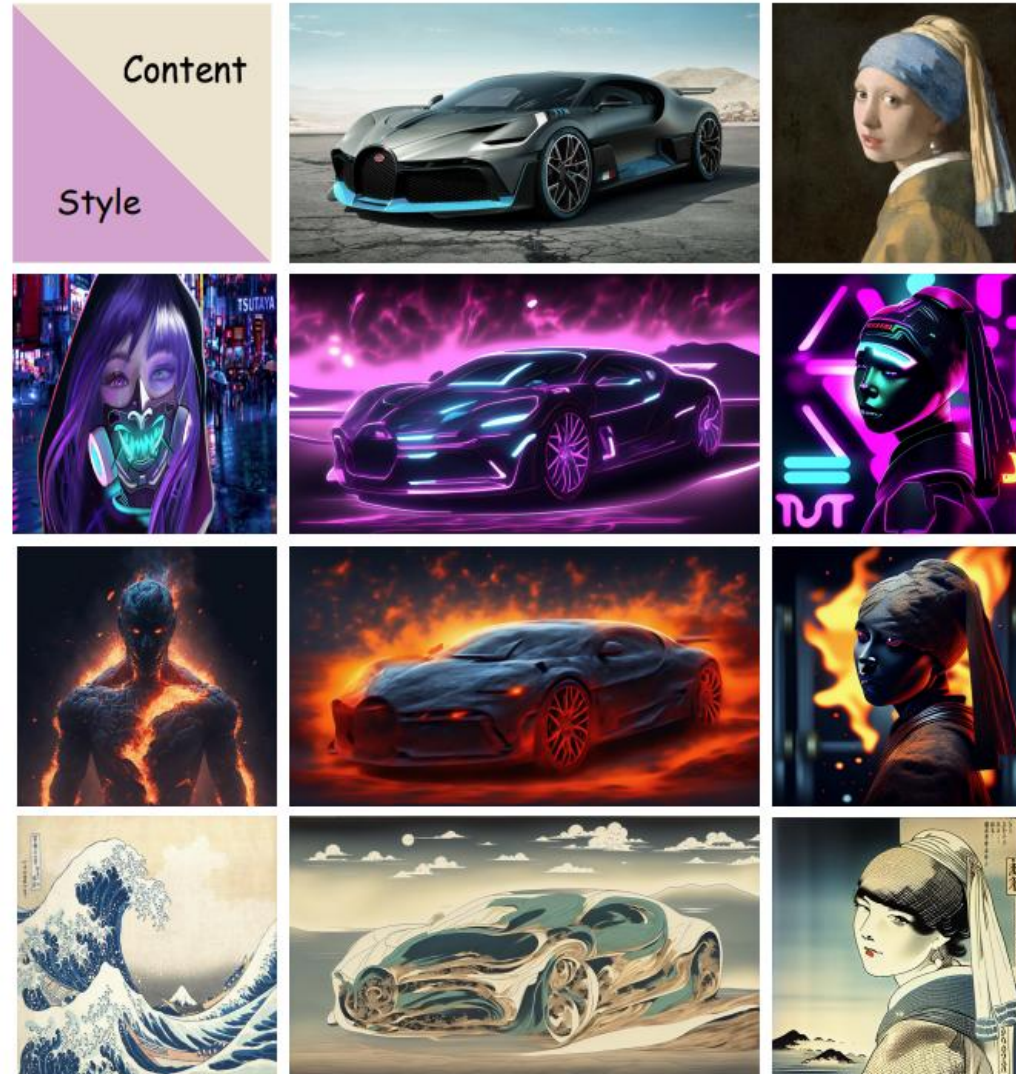
+STRE

+SERE

DEADiff

■ Experiments

Applications



Experiments

Ablations

Style 1



Style 2

"A motorcycle"

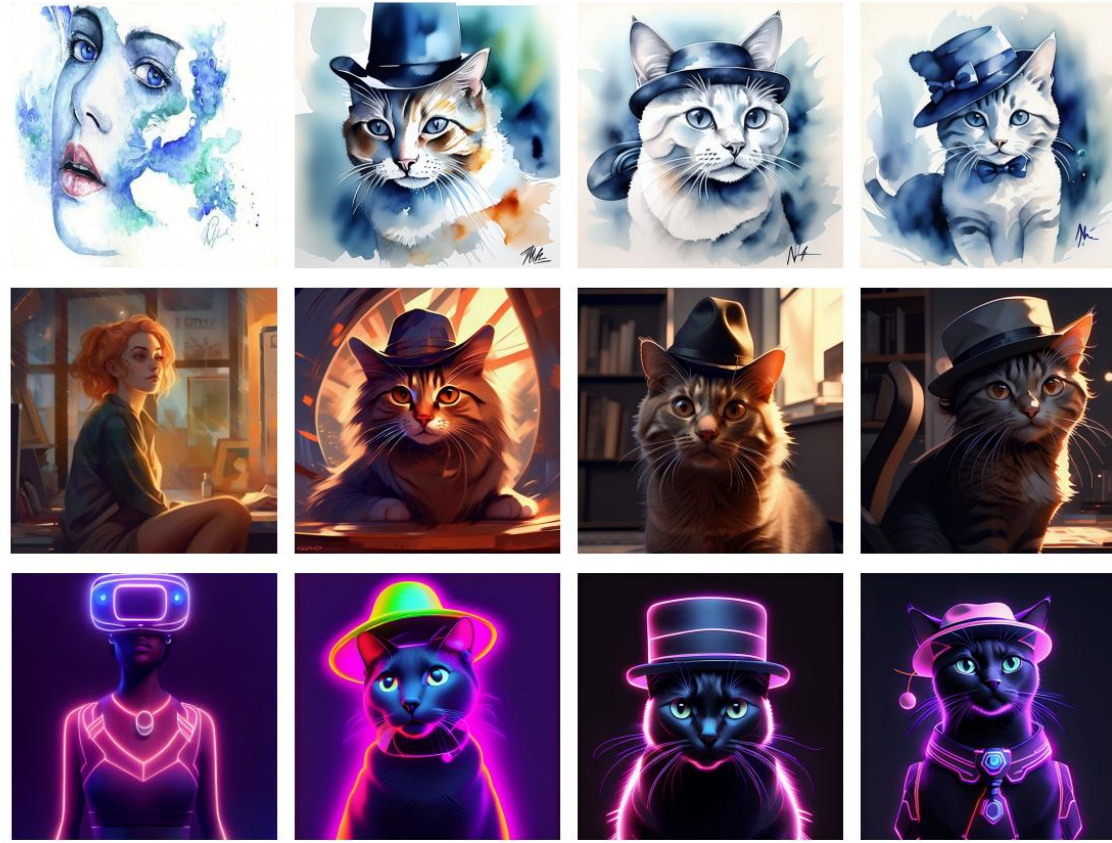


"A robot"

■ Experiments

Ablations

"A cat wearing a hat."



Reference

SD v1.5

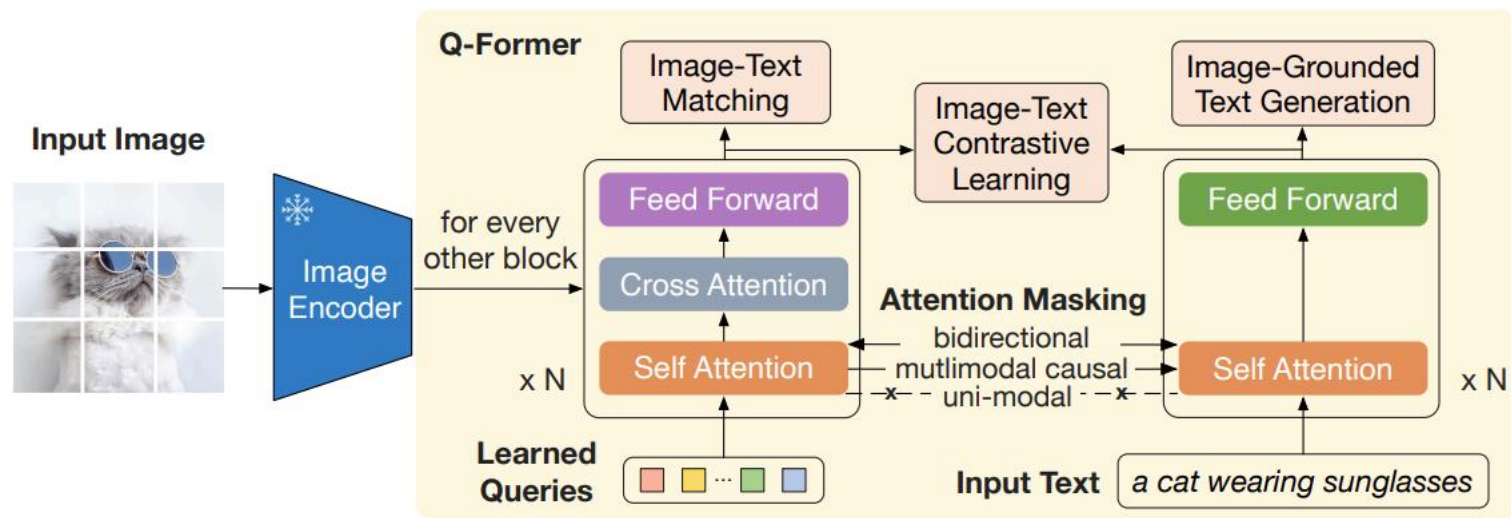
Realistic Vision V5.1

DreamShaper V8

Thanks!

Background

BLIP2 pre-training



Q: query token positions; **T**: text token positions.

■ masked □ unmasked

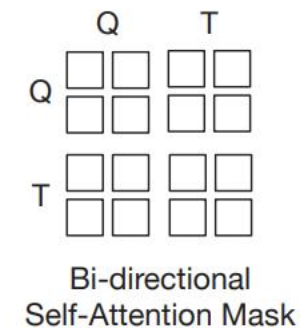


Image-Text Matching

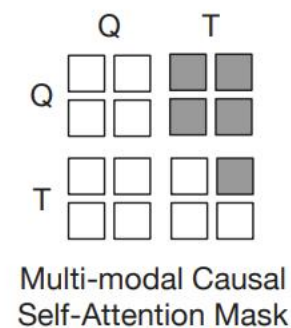


Image-Grounded Text Generation

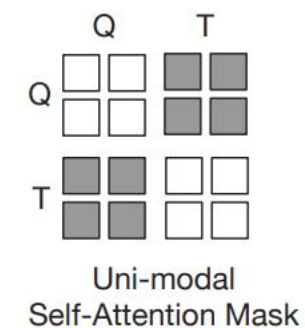


Image-Text Contrastive Learning