

MicroAST: Towards Super-Fast Ultra-Resolution Arbitrary Style Transfer

Zhizhong Wang, Lei Zhao*, Zhiwen Zuo, Ailin Li, Haibo Chen, Wei Xing*, Dongming Lu
AAAI 2023 Oral

STRUCT Group Seminar
Presenter: Zhengbo Xu
2023.01.29

OUTLINE

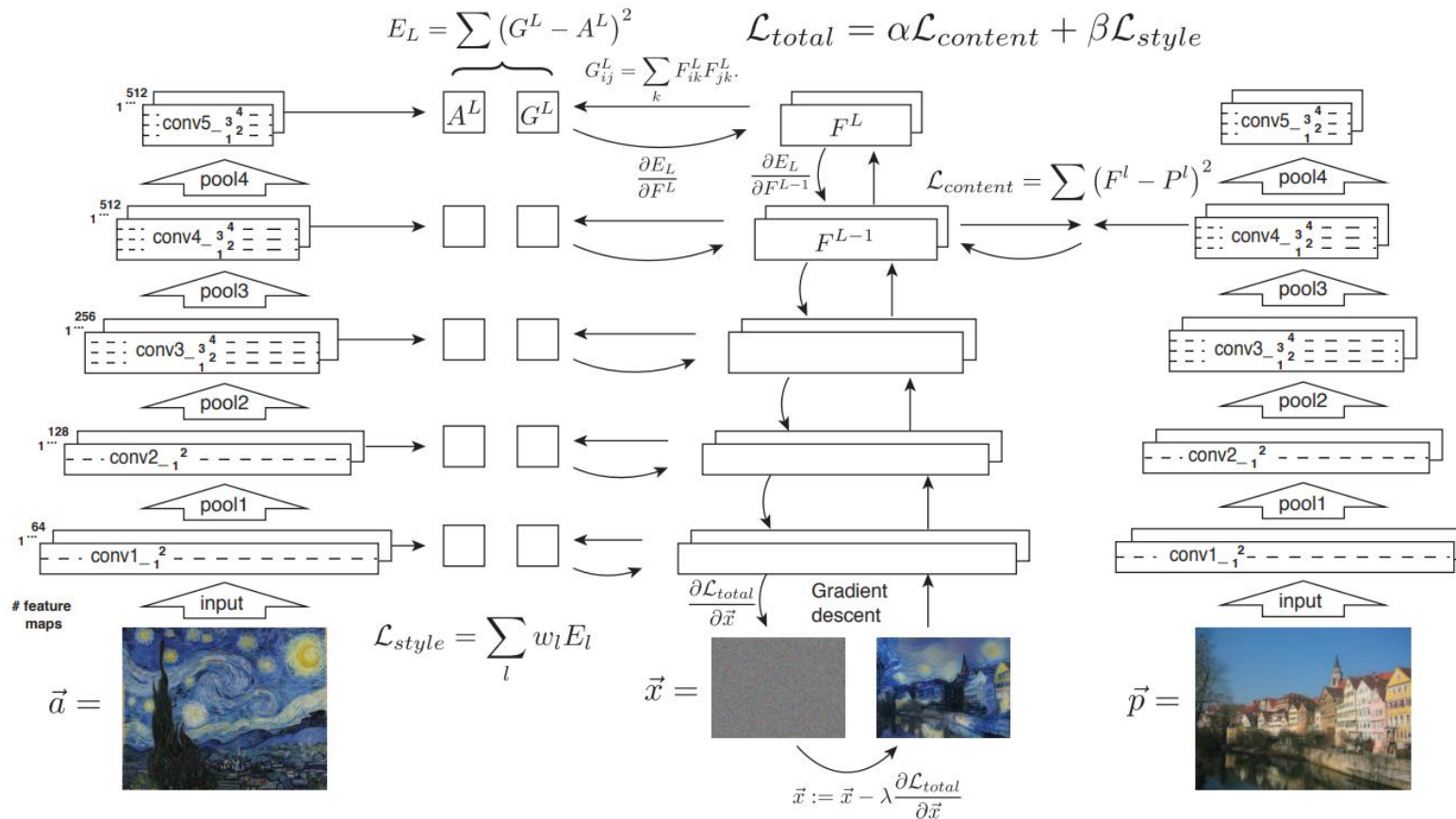
- Authorship
- Background
- Method
- Experiments
- Conclusion

OUTLINE

- Authorship
- **Background**
- Method
- Experiments
- Conclusion

BACKGROUND: Neural Style Transfer

Overview



BACKGROUND: Neural Style Transfer

Improved aspects

- Efficiency
- Quality
- Generalization
- Diversity
- User Control

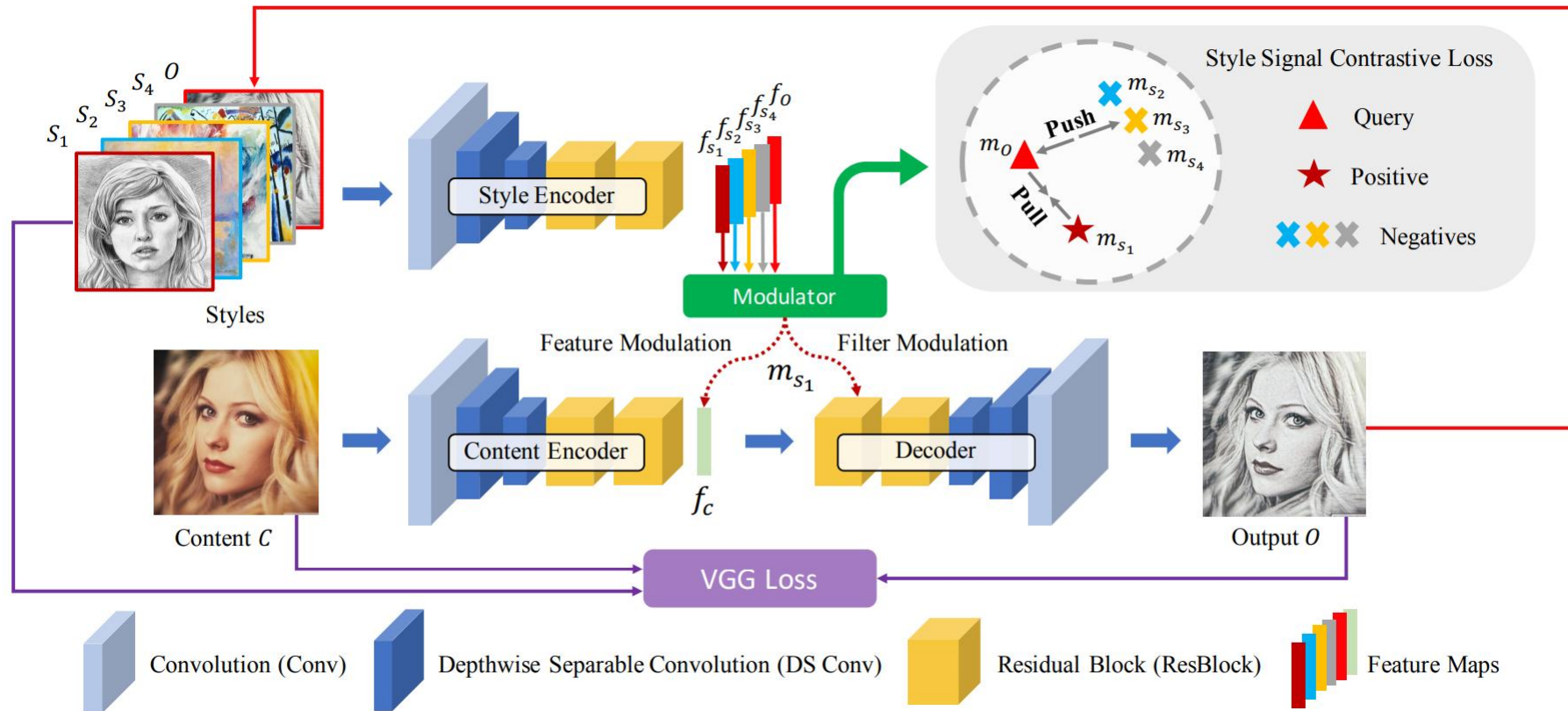
Generate 4K images with limited resources

OUTLINE

- Authorship
- Background
- Method
- Experiments
- Conclusion

METHOD

Overview



METHOD

Pipeline

- Extract features from content image C : $f_c := E_c(C)$.
- Extract features from style image S : $f_s := E_s(S)$.
- Convert f_s into style modulation signals: $m_s := \bar{\mathcal{M}}(f_s)$
- Stylize f_c using micro decoder D : $O := D(f_c, m_s)$

Training Loss

$$\mathcal{L}_{full} := \lambda_c \mathcal{L}_c + \lambda_s \mathcal{L}_s + \lambda_{ssc} \mathcal{L}_{ssc}$$

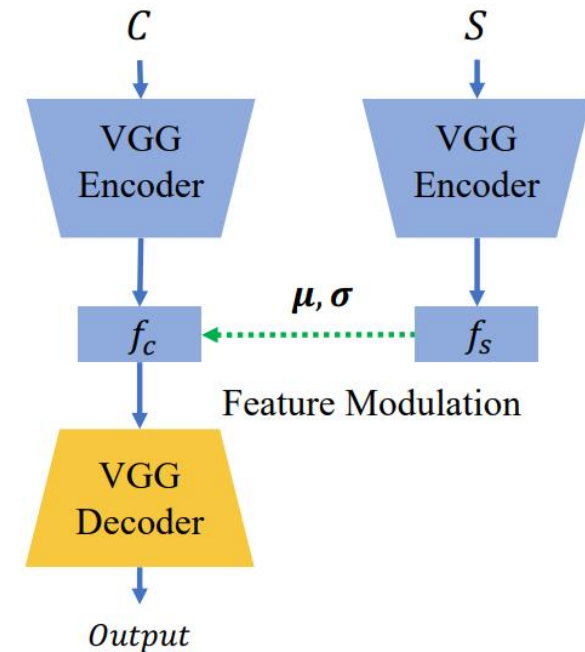
Modulation Strategies in AST

AdaIN (Huang and Belongie 2017)

$$\text{AdaIN}(f_c, f_s) := \sigma(f_s) \left(\frac{f_c - \mu(f_c)}{\sigma(f_c)} \right) + \mu(f_s).$$

Features & Requirements

- content and style encoder are identical
- encoders and decoder must be as complex as VGG
- encoders are fixed, decoder is trained



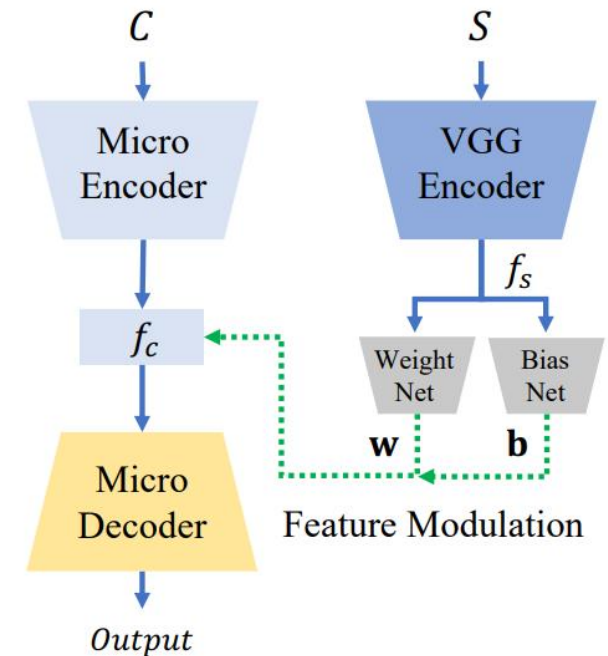
Modulation Strategies in AST

DIN (Jing et al. 2020)

$$DIN(f_c, f_s) := \mathbf{w} \left(\frac{f_c - \mu(f_c)}{\sigma(f_c)} \right) + \mathbf{b}.$$

Features & Requirements

- content encoder and decoder are lightweight
- style encoder is also high-cost VGG
- use subnets to predict weight and bias

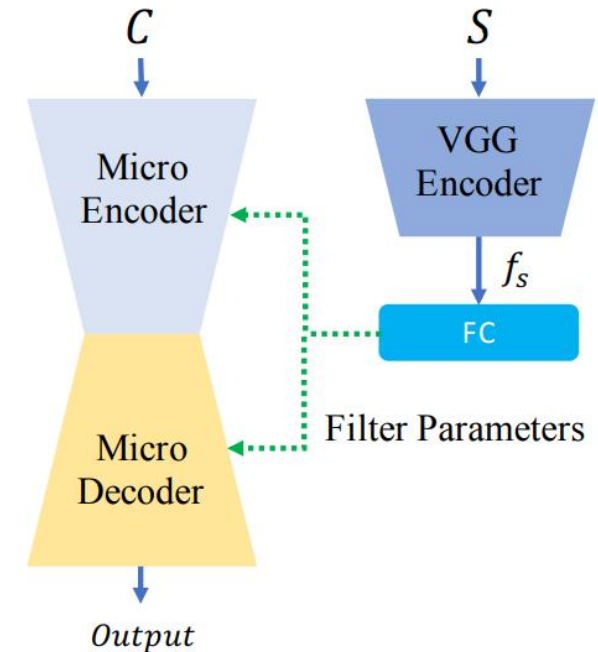


Modulation Strategies in AST

MetaNets (Shen, Yan, and Zeng 2018)

Features & Requirements

- content encoder and decoder are lightweight
- style encoder is also high-cost VGG
- style features help construct outputs
- FC layers lead to extra memory and slow inference time



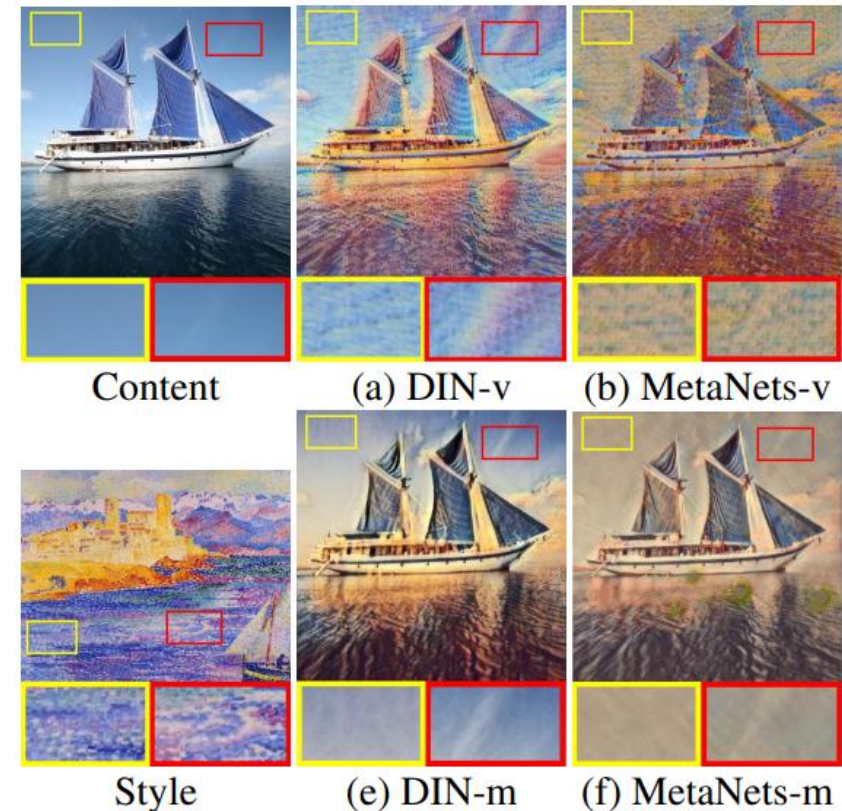
METHOD

Constraints

- The micro style encoder has limited ability to extract sufficiently complex style features
- The style signals are unitary and inflexible

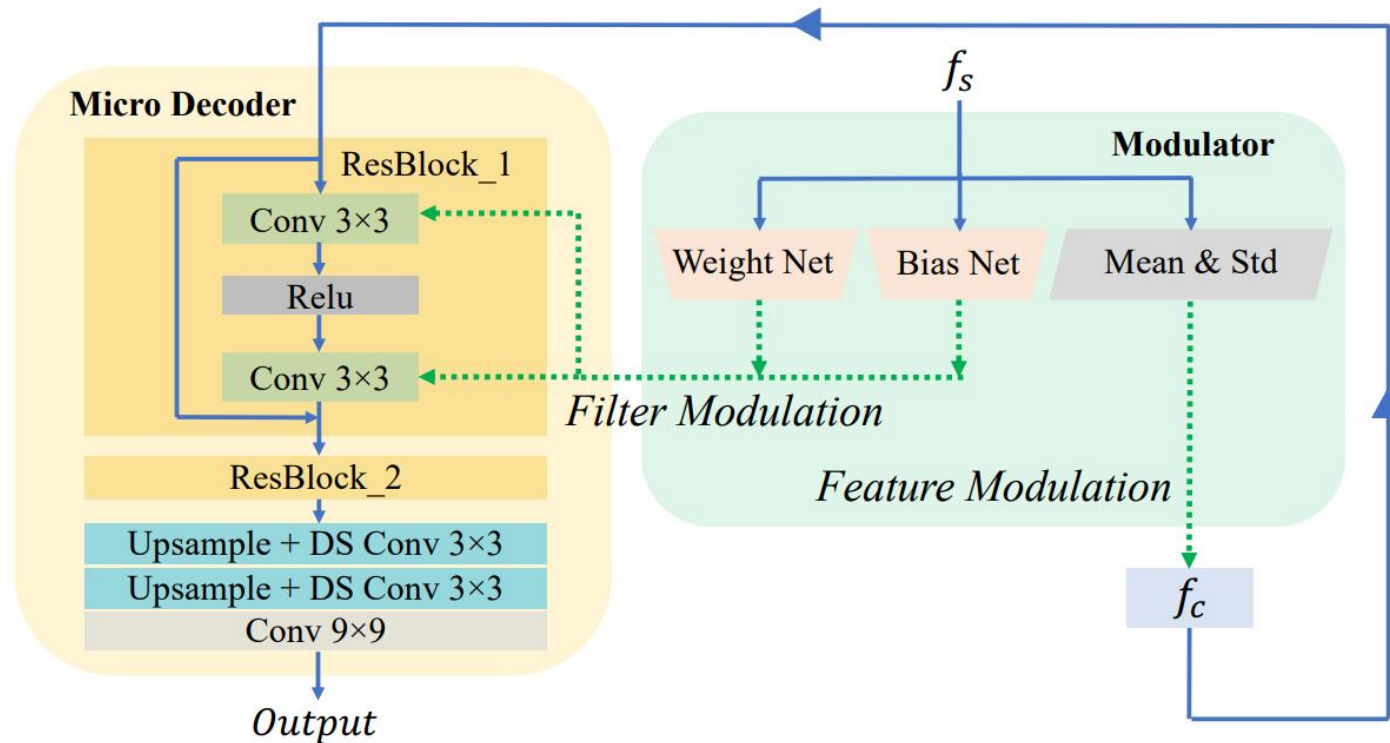
Solutions

- Introduce a dual-modulation strategy to inject more sophisticated and flexible style signals
- Propose a contrastive loss



METHOD

Dual-Modulation



$$m_s := (\boldsymbol{\mu}_s, \boldsymbol{\sigma}_s, \mathbf{w}_s, \mathbf{b}_s),$$

$$\text{DualMod}(D, f_c, m_s) := \text{FeatMod}(f_c, (\boldsymbol{\mu}_s, \boldsymbol{\sigma}_s)) + \text{FilterMod}(D, (\mathbf{w}_s, \mathbf{b}_s)),$$

METHOD

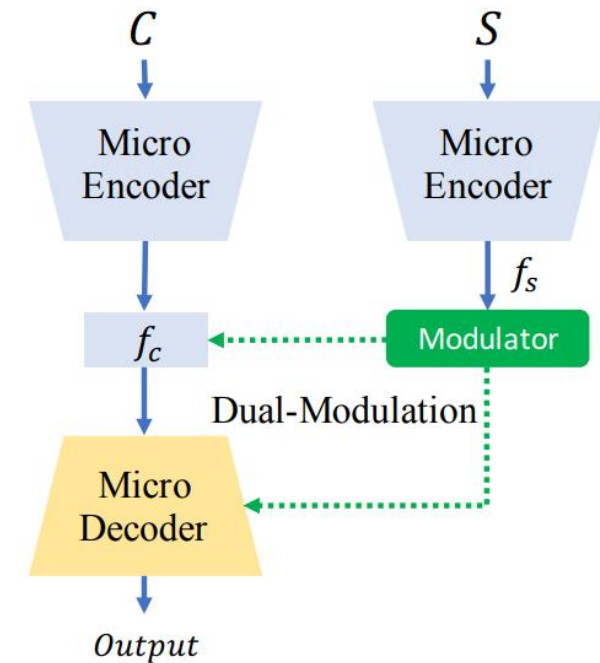
FeatMod

$$\boldsymbol{\mu}_s := \mu(f_s), \quad \boldsymbol{\sigma}_s := \sigma(f_s),$$

$$\text{FeatMod}(f_c, (\boldsymbol{\mu}_s, \boldsymbol{\sigma}_s)) := \boldsymbol{\sigma}_s \left(\frac{f_c - \mu(f_c)}{\sigma(f_c)} \right) + \boldsymbol{\mu}_s$$

Difference

- Generate from the micro encoder
- Just use mean and deviation



METHOD

FilterMod

$$\mathbf{w}_s := \xi_w(f_s), \quad \mathbf{b}_s := \xi_b(f_s),$$

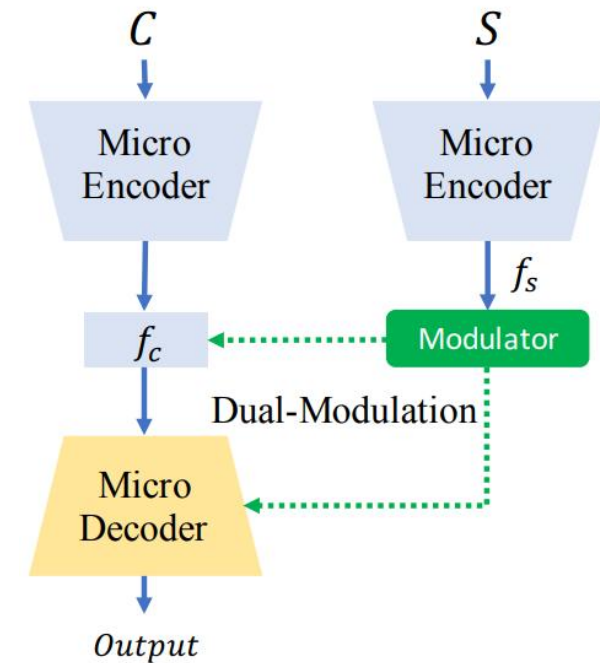
$$\text{FilterMod}(D, (\mathbf{w}_s, \mathbf{b}_s))$$

$$:= \text{ResBlock}(f_c, (\mathbf{w}_s, \mathbf{b}_s))$$

$$:= \text{Conv}(\text{Relu}(\text{Conv}(f_c, (\mathbf{w}_s, \mathbf{b}_s))), (\mathbf{w}_s, \mathbf{b}_s)) + f_c.$$

- Use simple subnets to predict weight and bias

$$\begin{aligned} \text{Conv}(f_c, (\mathbf{w}_s, \mathbf{b}_s)) &:= (\mathbf{w}_s * \mathcal{F} + \mathbf{b}_s) \circledast f_c \\ &:= (\mathbf{w}_s * \mathcal{F}) \circledast f_c + \mathbf{b}_s \circledast f_c \\ &:= \mathbf{w}_s * (\mathcal{F} \circledast f_c) + \mathbf{b}_s * f_c, \end{aligned}$$



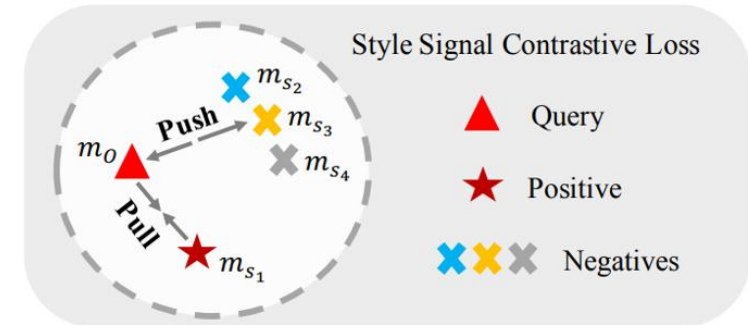
METHOD

Contrastive Learning

$$\mathcal{L}_{SSC} := \sum_{i=1}^N \frac{\|m_{o_i} - m_{s_i}\|_2}{\sum_{j \neq i}^N \|m_{o_i} - m_{s_j}\|_2}.$$

Difference





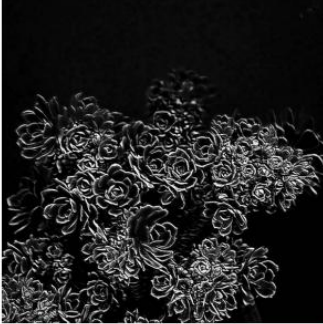






- Compare output images with style images
- A different form instead of vanilla loss (BCE)



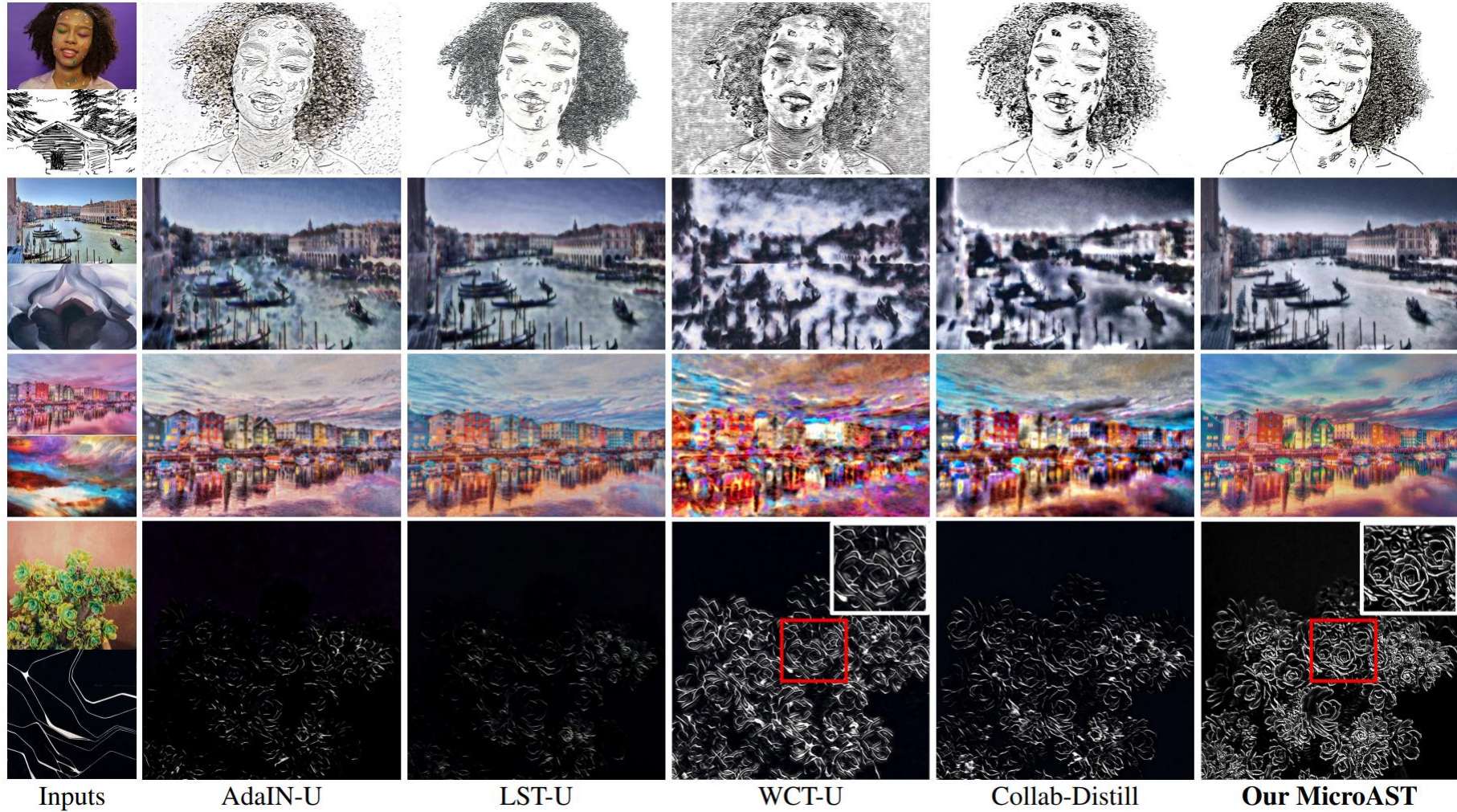
OUTLINE

- Authorship
- Background
- Method
- Experiments
- Conclusion

EXPERIMENTS

	Content Images		
			
Style Images	Stylized Outputs		
			
			

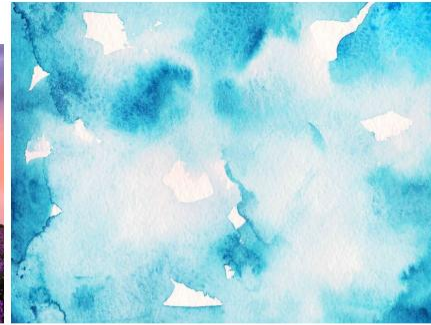
EXPERIMENTS



EXPERIMENTS



Content Image



Style Image



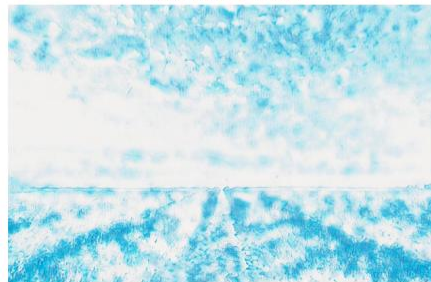
Our MicroAST



AdaIN [7]-U [2]



LST [10]-U [2]



WCT [11]-U [2]



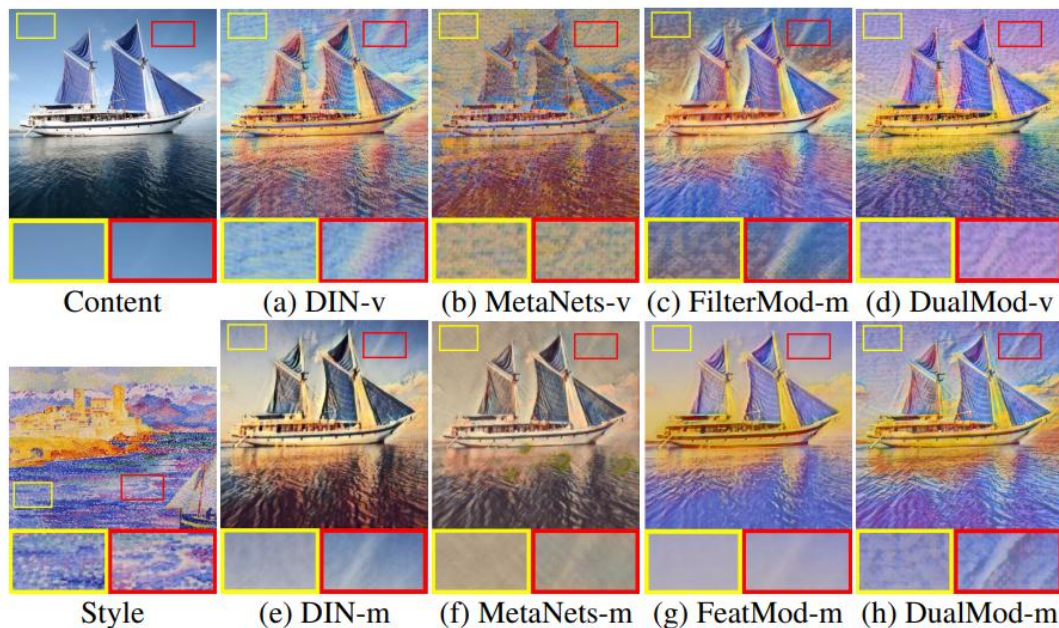
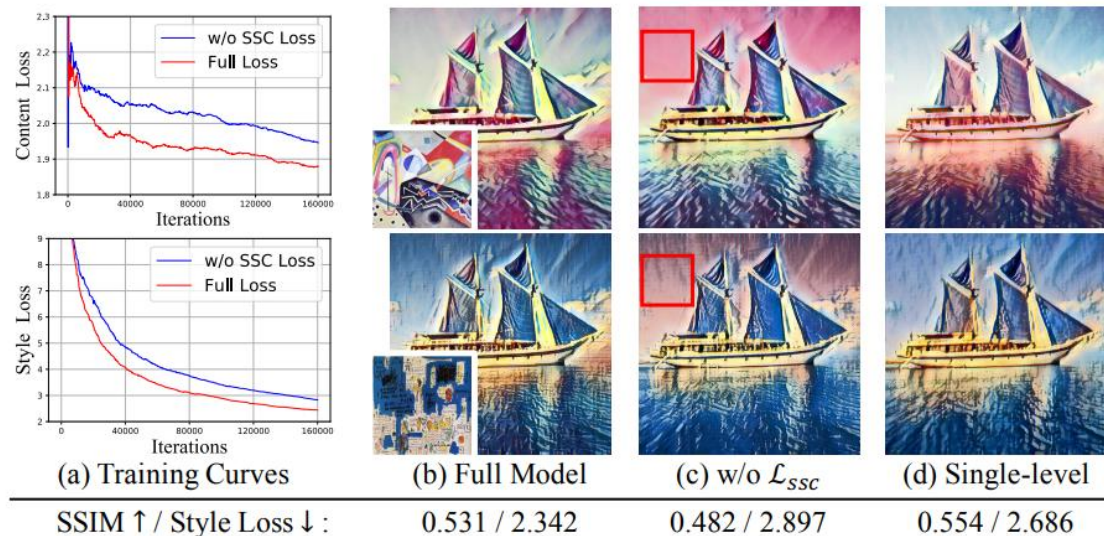
Collab-Distill [14]-U [2]

EXPERIMENTS



#Negative	1	7*	15	31
SSIM \uparrow	0.440	0.531	0.504	0.492
Style Loss \downarrow	2.547	2.342	2.333	2.332

EXPERIMENTS



OUTLINE

- Authorship
- Background
- Method
- Experiments
- Conclusion

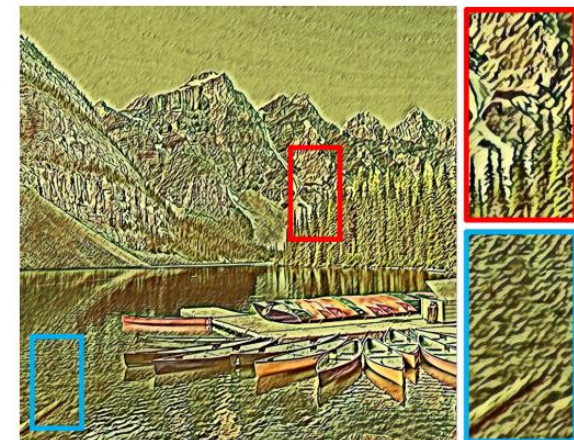
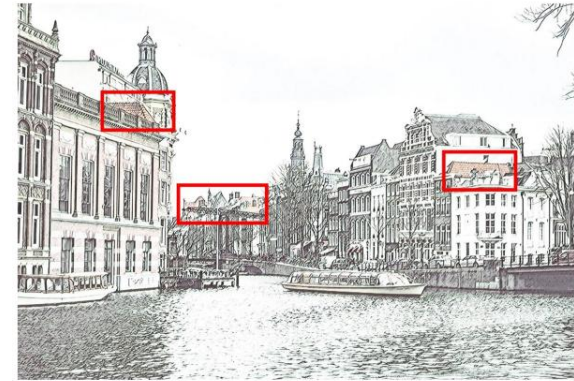
CONCLUSION

- Propose lightweight MicroAST to achieve super-fast ultra-resolution arbitrary style transfer
- Introduce the dual-modulation strategy
- Introduce a new style signal contrastive loss



LIMITATION

- Under-stylized results
- Fail to transfer complicated texture details
- Fail to deal with small stroke size



Thanks for listening!