

Sem-CS: Semantic CLIPStyler for Text-Based Image Style Transfer

Chanda Grover Kamra¹, Indra Deep Mastan², Debayan
Gupta¹

¹ Ashoka University, Computer Science, India.

² LNM Institute of Information Technology, Jaipur, India.

30th International Conference on Image Processing
ICIP 2023, Kuala Lumpur, Malaysia

Agenda

- Motivation
- Method
- Algorithm
- Results
- References

Motivation

- Over Stylization
 - *Distortion of content features.*

Motivation

- Over Stylization
 - *Distortion of content features.*

Style Text

A monet
style
painting

Input Image



CLIPStyler [1]



Gen-Art [2]



Sem-CS (Ours)



Motivation

- Over Stylization
 - *Distortion of content features.*
- Content Mismatch
 - *Style Spillover between dissimilar objects.*

Style Text

A monet
style
painting

Input Image



CLIPStyler [1]



Gen-Art [2]



Sem-CS (Ours)



Motivation

- Over Stylization
 - *Distortion of content features.*
- Content Mismatch
 - *Style Spillover between dissimilar objects.*

Style Text

A monet
style
painting

Input Image



CLIPStyler [1]



Gen-Art [2]



Sem-CS (Ours)



Style Text

Desert
Sand

Input Image



CLIPStyler [1]



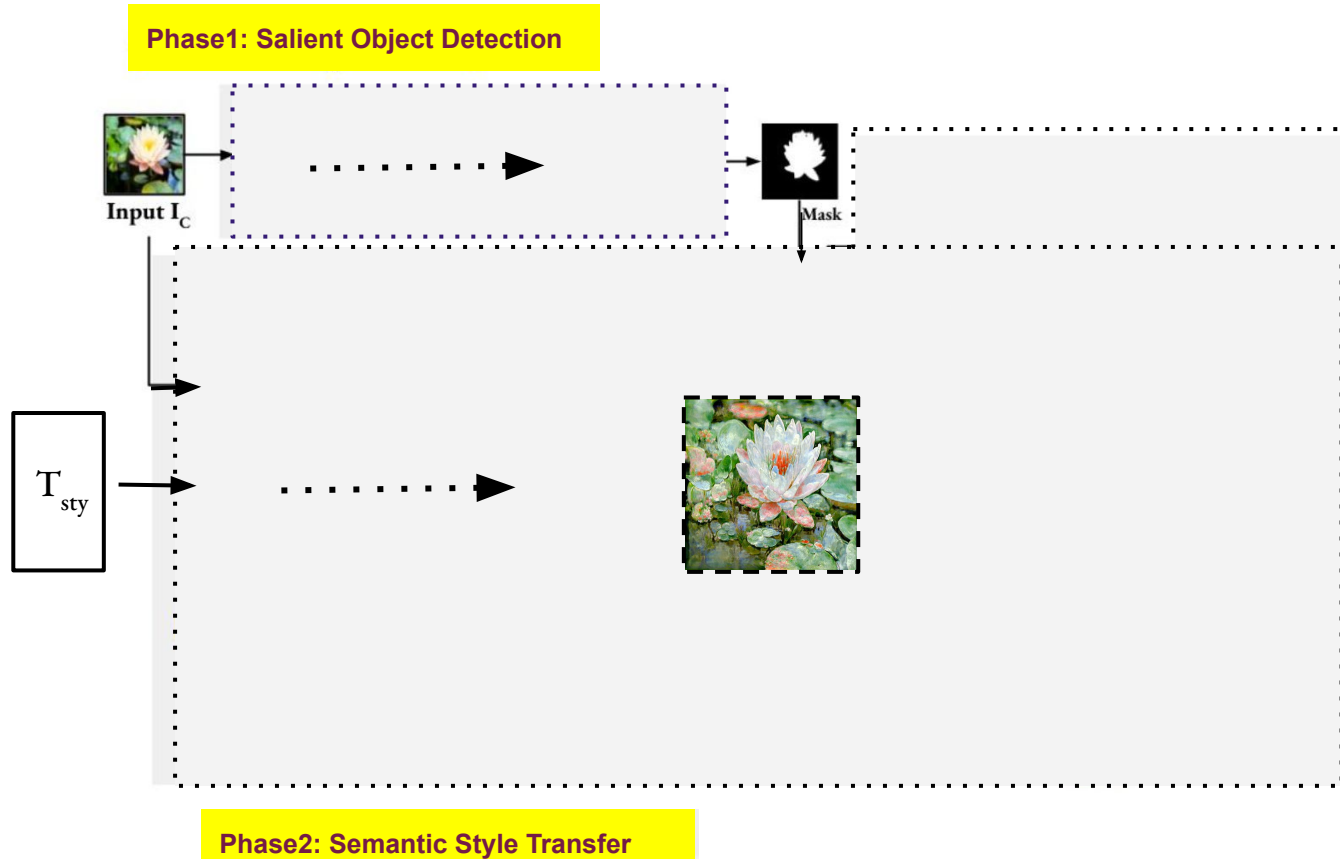
Gen-Art [2]



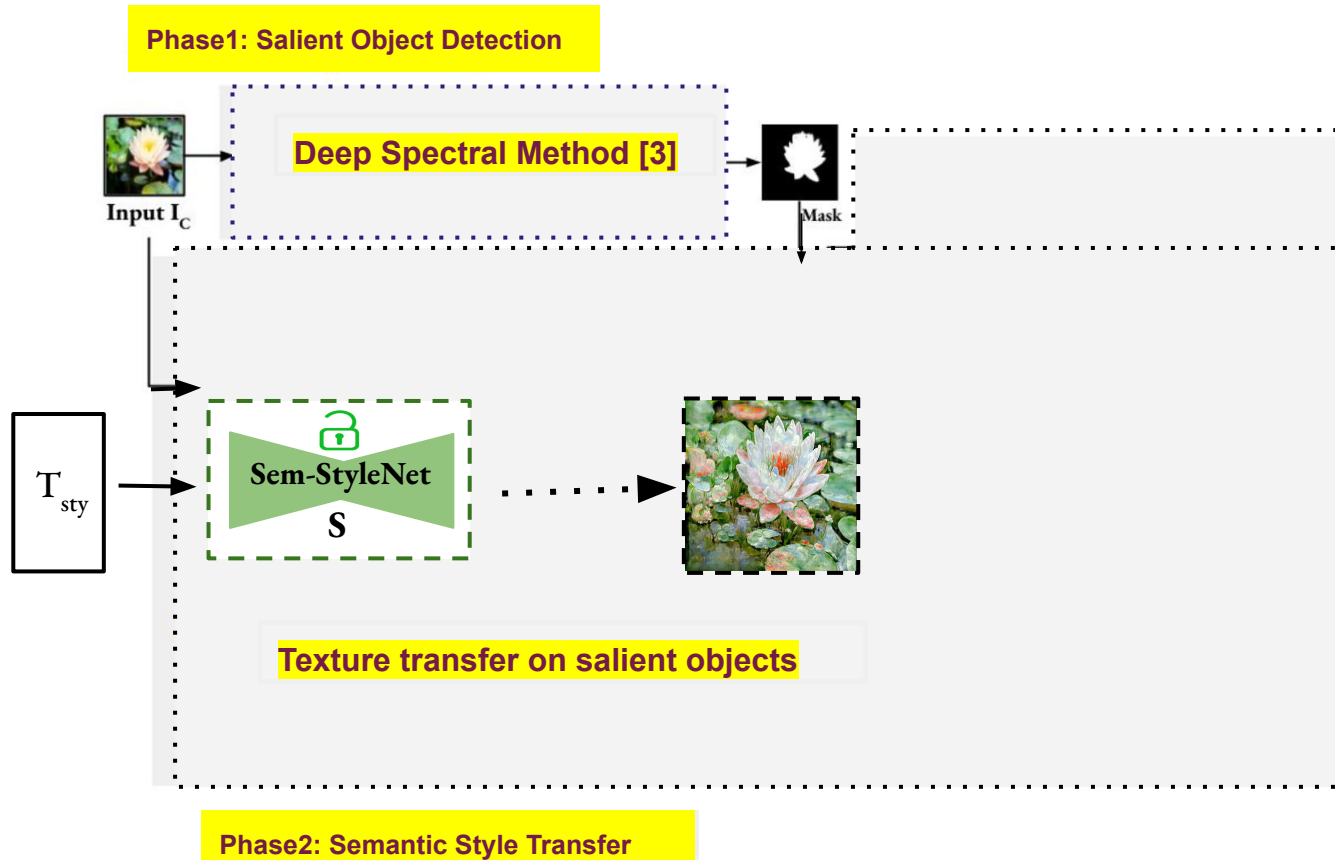
Sem-CS (Ours)



Method

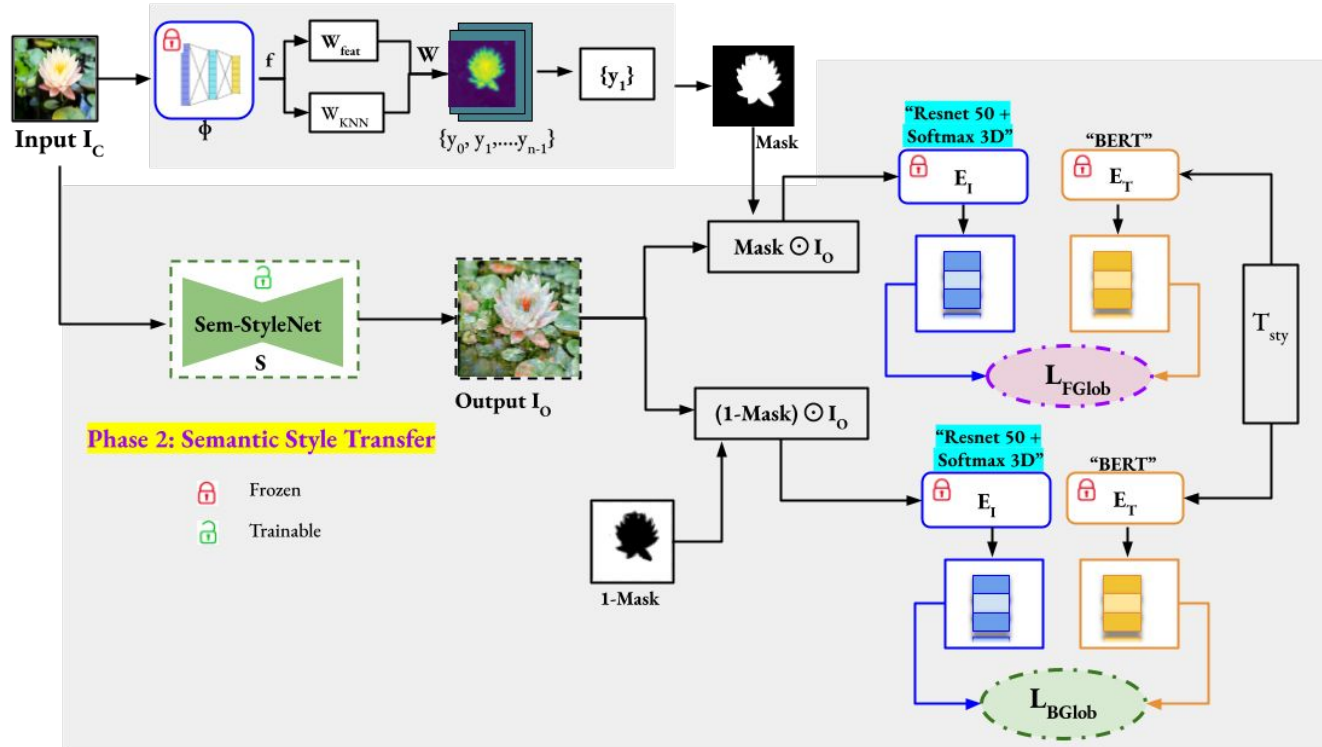


Method



Method: Sem-CS

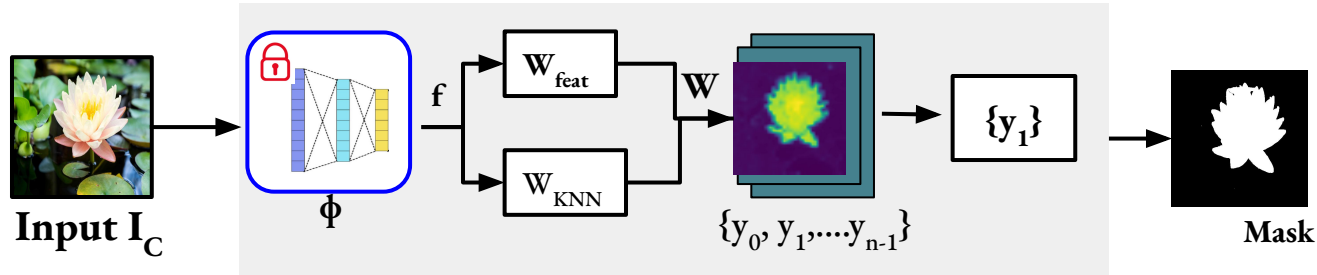
Phase 1: Salient Object Detection



Notations

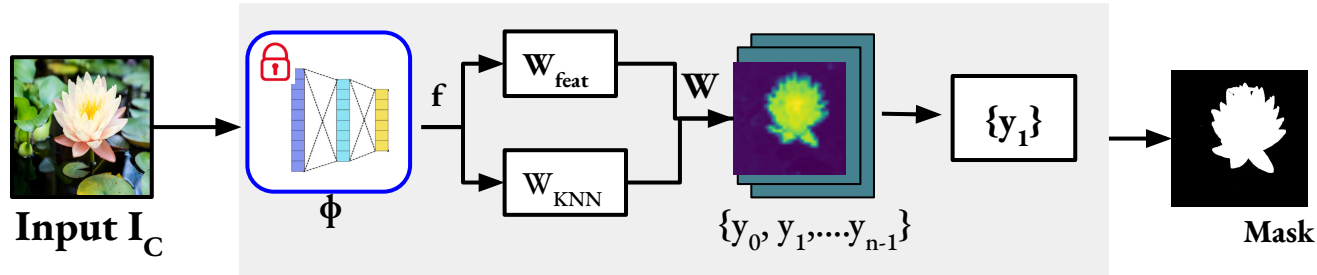
- I_C : Content Image
- T_{sty} : Style Text
- I_O : Stylized Output
- f : Deep patch features
- ϕ : Vision Transformer
- W_{KNN} : Color Matrix
- W_{feat} : Feature Matrix
- W : Semantic Affinity Matrix
- $\{y_0, y_1, \dots, y_{n-1}\}$: Eigen vectors
- S : Semantic StyleNet
- \odot : Hadamard Product
- E_I : CLIP Image Encoder
- E_T : CLIP Text Encoder
- L_{BGlob} : Global background loss
- L_{FGlob} : Global foreground loss

Method



- I_C : Content Image
- ϕ : Vision Transformer
- f : Pre-trained dense features
- W_{feat} : feature Matrix
- W_{KNN} : Colour Affinity Matrix
- W : Semantic Affinity Matrix
- $\{y_0, y_1, \dots, y_{n-1}\}$: Eigenvectors

Method

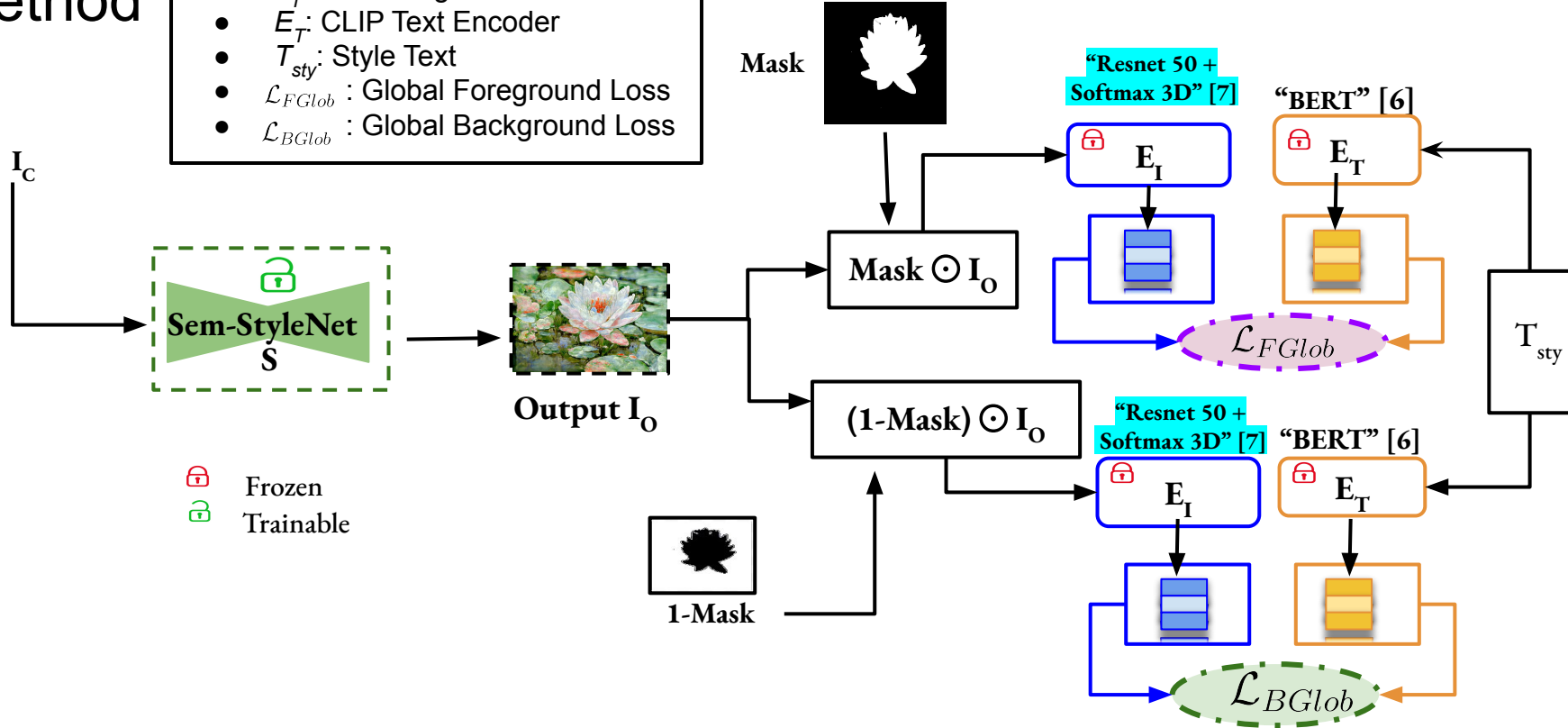


Phase1: Salient Object Detection

- I_C : Content Image
- ϕ : Vision Transformer
- f : Pre-trained dense features
- W_{feat} : feature Matrix
- W_{KNN} : Colour Affinity Matrix
- W : Semantic Affinity Matrix
- $\{y_0, y_1, \dots, y_{n-1}\}$: Eigenvectors

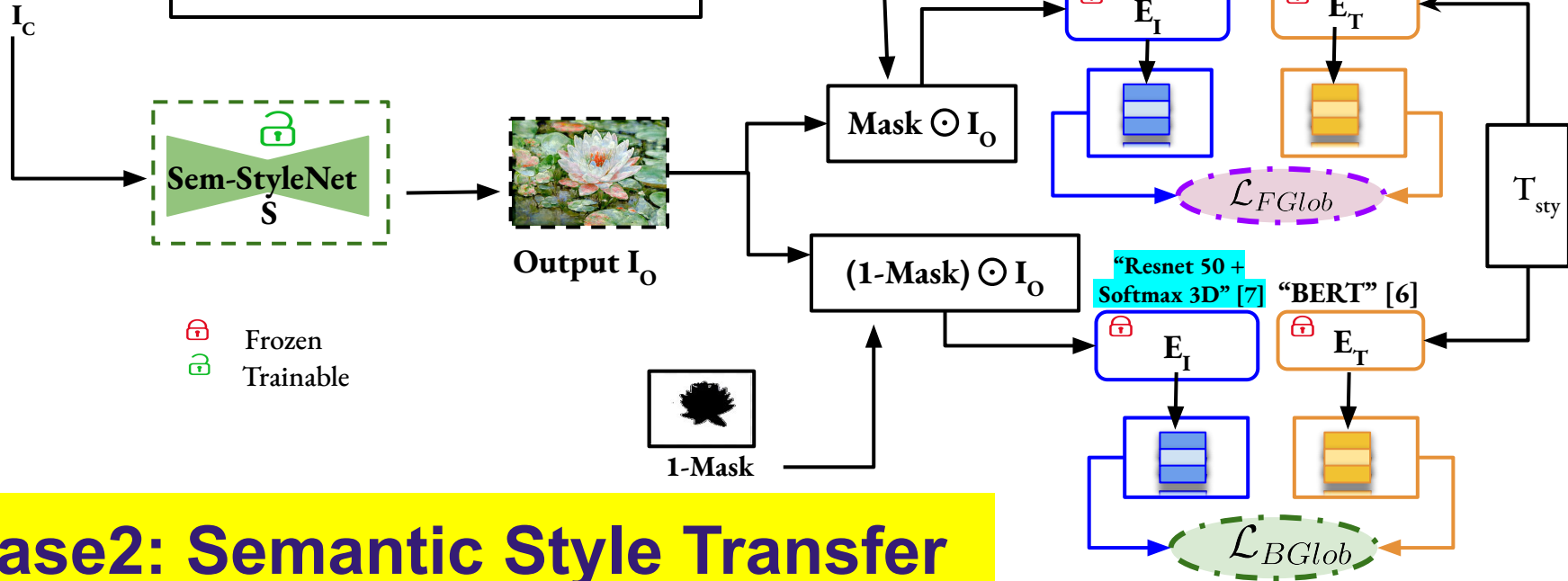
Method

- S : Semantic StyleNet
- I_o : Stylized Output
- E_i : CLIP Image Encoder
- E_T : CLIP Text Encoder
- T_{sty} : Style Text
- \mathcal{L}_{FGlob} : Global Foreground Loss
- \mathcal{L}_{BGlob} : Global Background Loss



Method

- S : Semantic StyleNet
- I_o : Stylized Output
- E_i : CLIP Image Encoder
- E_T : CLIP Text Encoder
- T_{sty} : Style Text
- \mathcal{L}_{FGlob} : Global Foreground Loss
- \mathcal{L}_{BGlob} : Global Background Loss



Phase2: Semantic Style Transfer

Algorithm: Semantic CLIPStyler

1: **SEM-CS**($I_C, \phi, T_{sty}, E_T, E_I, S$)

- I_C : Content Image
- ϕ : Vision Transformer
- T_{sty} : Style Text
- E_T : CLIP Text Encoder
- E_I : CLIP Image Encoder
- S : Semantic StyleNet

Algorithm: Semantic CLIPStyler

1: **SEM-CS**($I_C, \phi, T_{sty}, E_T, E_I, S$)

▷ *Compute Mask for salient objects identification*

2: $W = \text{AffinityMatrix}(I_C, \phi,)$

3: $\{y_0, y_1, \dots, y_{n-1}\} = \text{Eigen_Decomposition}(W)$

4: $Mask = \text{Extract_Salient_Object}(y_1)$

- I_C : Content Image
- ϕ : Vision Transformer
- T_{sty} : Style Text
- E_T : CLIP Text Encoder
- E_I : CLIP Image Encoder
- S : Semantic StyleNet

- $\{y_0, y_1, \dots, y_{n-1}\}$: Eigenvectors
- W : Semantic Affinity Matrix

Algorithm: Semantic CLIPStyler

- 1: **SEM-CS**($I_C, \phi, T_{sty}, E_T, E_I, S$)
 - ▷ *Compute Mask for salient objects identification*
- 2: $W = \text{AffinityMatrix}(I_C, \phi,)$
- 3: $\{y_0, y_1, \dots, y_{n-1}\} = \text{Eigen_Decomposition}(W)$
- 4: $Mask = \text{Extract_Salient_Object}(y_1)$

▷ *Perform Semantic Style Transfer*

- 5: $t_{fg}, t_{bg} = \text{Parse_Style_Text}(T_{sty})$
- 6: $I_{fg}, I_{bg} = Mask \odot S(I_C), (1 - Mask) \odot S(I_C)$

- I_C : Content Image
- ϕ : Vision Transformer
- T_{sty} : Style Text
- E_T : CLIP Text Encoder
- E_I : CLIP Image Encoder
- S : Semantic StyleNet
- $\{y_0, y_1, \dots, y_{n-1}\}$: Eigenvectors
- W : Semantic Affinity Matrix

- t_{fg} : Foreground Text embeddings
- t_{bg} : Background Text embeddings
- I_{fg} : Foreground Image embeddings
- I_{bg} : Background Image embeddings
- \odot : Hadamard Product

Algorithm: Semantic CLIPStyler

- 1: **SEM-CS**($I_C, \phi, T_{sty}, E_T, E_I, S$)
 - ▷ *Compute Mask for salient objects identification*
- 2: $W = \text{AffinityMatrix}(I_C, \phi,)$
- 3: $\{y_0, y_1, \dots, y_{n-1}\} = \text{Eigen_Decomposition}(W)$
- 4: $Mask = \text{Extract_Salient_Object}(y_1)$
 - ▷ *Perform Semantic Style Transfer*
- 5: $t_{fg}, t_{bg} = \text{Parse_Style_Text}(T_{sty})$
- 6: $I_{fg}, I_{bg} = Mask \odot S(I_C), (1 - Mask) \odot S(I_C)$

- I_C : Content Image
- ϕ : Vision Transformer
- T_{sty} : Style Text
- E_T : CLIP Text Encoder
- E_I : CLIP Image Encoder
- S : Semantic StyleNet
- $\{y_0, y_1, \dots, y_{n-1}\}$: Eigenvectors
- W : Semantic Affinity Matrix

- t_{fg} : Foreground Text embeddings
- t_{bg} : Background Text embeddings
- I_{fg} : Foreground Image embeddings
- I_{bg} : Background Image embeddings
- \odot : Hadamard Product

- \mathcal{L}_{FGlob} : Global Foreground Loss

▷ *Global Foreground Loss*

- 7: Compute Foreground Image Direction Loss $\Delta f g_I = E_I(I_{fg}) - E_I(I_C)$
- 8: Compute Foreground Text Direction Loss $\Delta f g_T = E_T(t_{fg}) - E_T(t_{src})$
- 9: $\mathcal{L}_{FGlob} = \text{Cosine_similarity}(\Delta f g_I, \Delta f g_T) = 1 - \frac{\Delta f g_I \cdot \Delta f g_T}{|\Delta f g_I| |\Delta f g_T|}$

Algorithm: Semantic CLIPStyler

1: **SEM-CS**($I_C, \phi, T_{sty}, E_T, E_I, S$)

▷ *Compute Mask for salient objects identification*

2: $W = \text{AffinityMatrix}(I_C, \phi,)$

3: $\{y_0, y_1, \dots, y_{n-1}\} = \text{Eigen_Decomposition}(W)$

4: $Mask = \text{Extract_Salient_Object}(y_1)$

▷ *Perform Semantic Style Transfer*

5: $t_{fg}, t_{bg} = \text{Parse_Style_Text}(T_{sty})$

6: $I_{fg}, I_{bg} = Mask \odot S(I_C), (1 - Mask) \odot S(I_C)$

▷ *Global Foreground Loss*

7: Compute Foreground Image Direction Loss $\Delta fg_I = E_I(I_{fg}) - E_I(I_C)$

8: Compute Foreground Text Direction Loss $\Delta fg_T = E_T(t_{fg}) - E_T(t_{src})$

9: $\mathcal{L}_{FGlob} = \text{Cosine_similarity}(\Delta fg_I, \Delta fg_T) = 1 - \frac{\Delta fg_I \cdot \Delta fg_T}{|\Delta fg_I| |\Delta fg_T|}$

▷ *Global Background Loss*

10: Compute Background Image Direction Loss $\Delta bg_I = E_I(I_{bg}) - E_I(I_C)$

11: Compute Background Text Direction Loss $\Delta bg_T = E_T(t_{bg}) - E_T(t_{src})$

12: $\mathcal{L}_{BGlob} = \text{Cosine_similarity}(\Delta bg_I, \Delta bg_T) = 1 - \frac{\Delta bg_I \cdot \Delta bg_T}{|\Delta bg_I| |\Delta bg_T|}$

- I_C : Content Image
- ϕ : Vision Transformer
- T_{sty} : Style Text
- E_T : CLIP Text Encoder
- E_I : CLIP Image Encoder
- S : Semantic StyleNet
- $\{y_0, y_1, \dots, y_{n-1}\}$: Eigenvectors
- W : Semantic Affinity Matrix
- t_{fg} : Foreground Text embeddings
- t_{bg} : Background Text embeddings
- I_{fg} : Foreground Image embeddings
- I_{bg} : Background Image embeddings
- \odot : Hadamard Product
- \mathcal{L}_{FGlob} : Global Foreground Loss
- \mathcal{L}_{BGlob} : Global Background Loss

Algorithm: Semantic CLIPStyler

- 1: **SEM-CS**($I_C, \phi, T_{sty}, E_T, E_I, S$)
 - ▷ *Compute Mask for salient objects identification*
- 2: $W = \text{AffinityMatrix}(I_C, \phi,)$
- 3: $\{y_0, y_1, \dots, y_{n-1}\} = \text{Eigen_Decomposition}(W)$
- 4: $Mask = \text{Extract_Salient_Object}(y_1)$
- ▷ *Perform Semantic Style Transfer*
- 5: $t_{fg}, t_{bg} = \text{Parse_Style_Text}(T_{sty})$
- 6: $I_{fg}, I_{bg} = Mask \odot S(I_C), (1 - Mask) \odot S(I_C)$
- ▷ *Global Foreground Loss*
- 7: Compute Foreground Image Direction Loss $\Delta fg_I = E_I(I_{fg}) - E_I(I_C)$
- 8: Compute Foreground Text Direction Loss $\Delta fg_T = E_T(t_{fg}) - E_T(t_{src})$
- 9: $\mathcal{L}_{FGlob} = \text{Cosine_similarity}(\Delta fg_I, \Delta fg_T) = 1 - \frac{\Delta fg_I \cdot \Delta fg_T}{|\Delta fg_I| |\Delta fg_T|}$
- ▷ *Global Background Loss*
- 10: Compute Background Image Direction Loss $\Delta bg_I = E_I(I_{bg}) - E_I(I_C)$
- 11: Compute Background Text Direction Loss $\Delta bg_T = E_T(t_{bg}) - E_T(t_{src})$
- 12: $\mathcal{L}_{BGlob} = \text{Cosine_similarity}(\Delta bg_I, \Delta bg_T) = 1 - \frac{\Delta bg_I \cdot \Delta bg_T}{|\Delta bg_I| |\Delta bg_T|}$
- ▷ *Minimize loss and compute output I_O*
- 13: $I_O = \min_{\theta_S} (\mathcal{L}_{FGlob} + \lambda_{bg} \mathcal{L}_{BGlob})$

















- I_C : Content Image
- ϕ : Vision Transformer
- T_{sty} : Style Text
- E_T : CLIP Text Encoder
- E_I : CLIP Image Encoder
- S : Semantic StyleNet
- $\{y_0, y_1, \dots, y_{n-1}\}$: Eigenvectors
- W : Semantic Affinity Matrix
- t_{fg} : Foreground Text embeddings
- t_{bg} : Background Text embeddings
- I_{fg} : Foreground Image embeddings
- I_{bg} : Background Image embeddings
- \odot : Hadamard Product
- \mathcal{L}_{FGlob} : Global Foreground Loss
- \mathcal{L}_{BGlob} : Global Background Loss

Algorithm: Semantic CLIPStyler

















- 1: **SEM-CS**($I_C, \phi, T_{sty}, E_T, E_I, S$)
 - ▷ *Compute Mask for salient objects identification*
- 2: $W = \text{AffinityMatrix}(I_C, \phi,)$
- 3: $\{y_0, y_1, \dots, y_{n-1}\} = \text{Eigen_Decomposition}(W)$
- 4: $Mask = \text{Extract_Salient_Object}(y_1)$
 - ▷ *Perform Semantic Style Transfer*
- 5: $t_{fg}, t_{bg} = \text{Parse_Style_Text}(T_{sty})$
- 6: $I_{fg}, I_{bg} = Mask \odot S(I_C), (1 - Mask) \odot S(I_C)$
 - ▷ *Global Foreground Loss*
- 7: Compute Foreground Image Direction Loss $\Delta fg_I = E_I(I_{fg}) - E_I(I_C)$
- 8: Compute Foreground Text Direction Loss $\Delta fg_T = E_T(t_{fg}) - E_T(t_{src})$
- 9: $\mathcal{L}_{FGlob} = \text{Cosine_similarity}(\Delta fg_I, \Delta fg_T) = 1 - \frac{\Delta fg_I \cdot \Delta fg_T}{|\Delta fg_I| |\Delta fg_T|}$
 - ▷ *Global Background Loss*
- 10: Compute Background Image Direction Loss $\Delta bg_I = E_I(I_{bg}) - E_I(I_C)$
- 11: Compute Background Text Direction Loss $\Delta bg_T = E_T(t_{bg}) - E_T(t_{src})$
- 12: $\mathcal{L}_{BGlob} = \text{Cosine_similarity}(\Delta bg_I, \Delta bg_T) = 1 - \frac{\Delta bg_I \cdot \Delta bg_T}{|\Delta bg_I| |\Delta bg_T|}$
 - ▷ *Minimize loss and compute output I_O*
- 13: $I_O = \min_{\theta_S} (\mathcal{L}_{FGlob} + \lambda_{bg} \mathcal{L}_{BGlob})$

- I_C : Content Image
- ϕ : Vision Transformer
- T_{sty} : Style Text
- E_T : CLIP Text Encoder
- E_I : CLIP Image Encoder
- S : Semantic StyleNet
- $\{y_0, y_1, \dots, y_{n-1}\}$: Eigenvectors
- W : Semantic Affinity Matrix
- t_{fg} : Foreground Text embeddings
- t_{bg} : Background Text embeddings
- I_{fg} : Foreground Image embeddings
- I_{bg} : Background Image embeddings
- \odot : Hadamard Product
- \mathcal{L}_{FGlob} : Global Foreground Loss
- \mathcal{L}_{BGlob} : Global Background Loss

Results: Single Text Condition

Input Image	Style Text	CLIPStyler [1]	Gen-Art [2]	Sem-CS (Ours)
	Acrylic painting			
	A graffiti style painting			
	A fauvism style painting			
	An oil painting of white roses			

Results: Single Text Condition

Input Image	Style Text	CLIPStyler [1]	Gen-Art [2]	Sem-CS (Ours)
	Snowy			
	Red rocks			
	A watercolor painting with purple brush			
	A watercolor painting of leaf			

Results: Single Text Condition

Scores	CLIPStyler [1]	Gen-Art [2]	Sem-CS (Ours)
DISTS↑ [8]	0.32	0.25	0.34
NIMA↑ [9]	4.61	4.34	5.34
USer Study↑	28.3	33.1	38.4

Results: Double Text Condition

Style Text

F: Red Rocks

B: Snowy

Input Image



Gen-Art [2]

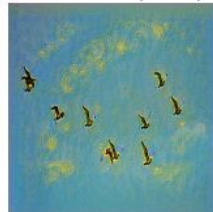
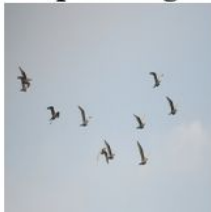


Sem-CS (ours)



F: Pop Art

B: Starry Night
By Vincent Van
Gogh



Results: Double Text Condition

Style Text

Input Image Gen-Art [2] Sem-CS (ours)



F: Red Rocks
B: Snowy

F: Pop Art
B: Starry Night
By Vincent Van Gogh



Scores	Gen-Art [2]	Sem-CS (Ours)
DISTS↑ [8]	0.20	0.33
NIMA↑ [9]	4.24	5.52
User Study ↑	48.2	51.7

References [I]

1. Kwon, Gihyun, and Jong Chul Ye. "Clipstyler: Image style transfer with a single text condition." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.
2. Yang, Zhenling, Huacheng Song, and Qiunan Wu. "Generative Artisan: A Semantic-Aware and Controllable CLIPstyler." arXiv preprint arXiv:2207.11598 (2022).
3. Melas-Kyriazi, Luke, et al. "Deep spectral methods: A surprisingly strong baseline for unsupervised semantic segmentation and localization." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.

References [II]

4. Long, Jonathan, Evan Shelhamer, and Trevor Darrell. "Fully convolutional networks for semantic segmentation." Proceedings of the IEEE conference on computer vision and pattern recognition. 2015.
5. Radford, Alec, et al. "Learning transferable visual models from natural language supervision." International conference on machine learning. PMLR, 2021.
6. Kenton, Jacob Devlin Ming-Wei Chang, and Lee Kristina Toutanova. "Bert: Pre-training of deep bidirectional transformers for language understanding." Proceedings of naacL-HLT. Vol. 1. 2019.
7. Wang, Pei, Yijun Li, and Nuno Vasconcelos. "Rethinking and improving the robustness of image style transfer." Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2021.

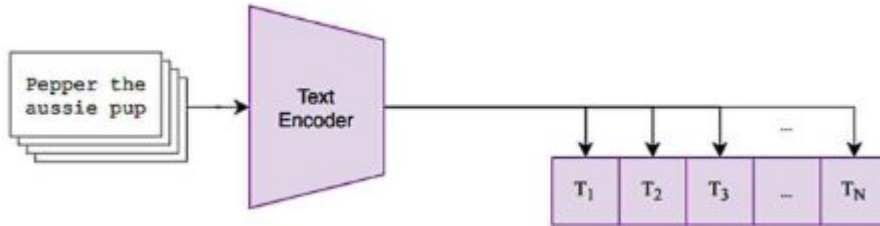
References [III]

8. Ding, Keyan, et al. "Image quality assessment: Unifying structure and texture similarity." *IEEE transactions on pattern analysis and machine intelligence* 44.5 (2020): 2567-2581.
9. Talebi, Hossein, and Peyman Milanfar. "NIMA: Neural image assessment." *IEEE transactions on image processing* 27.8 (2018): 3998-4011.

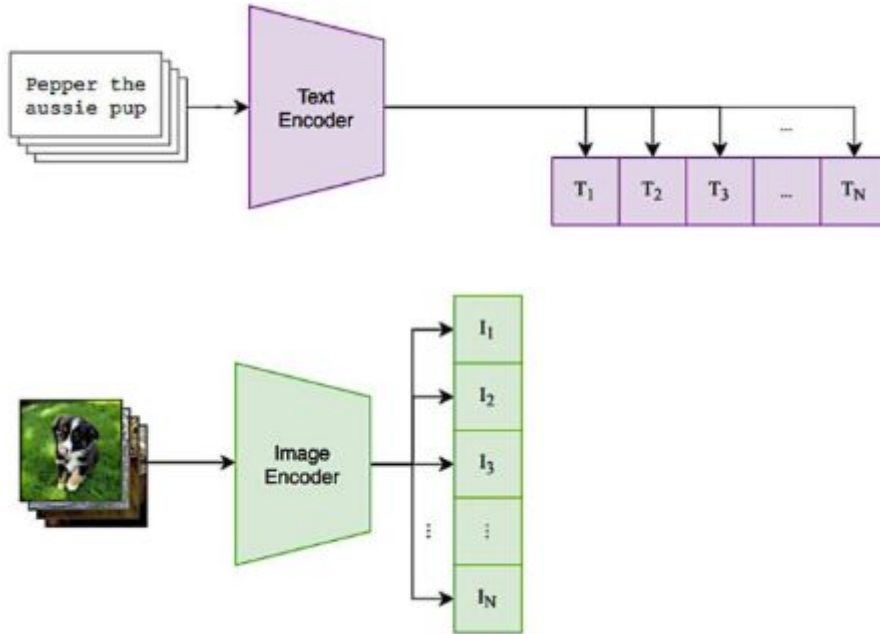
Thank You

Queries?

Background

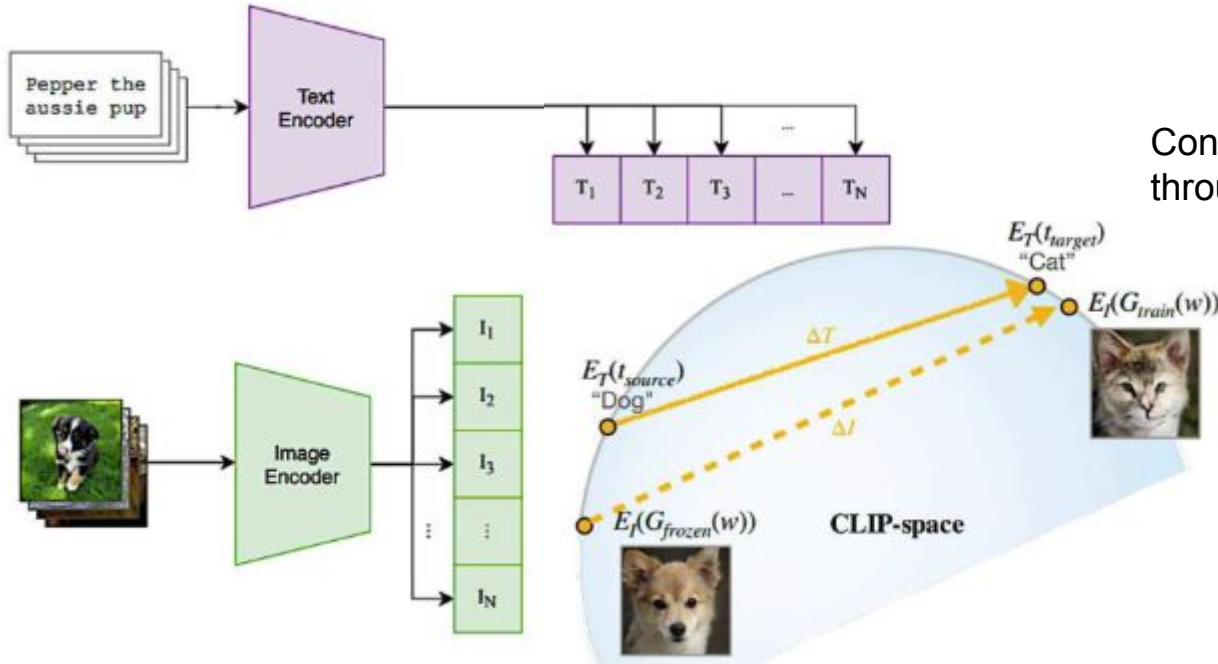


Background



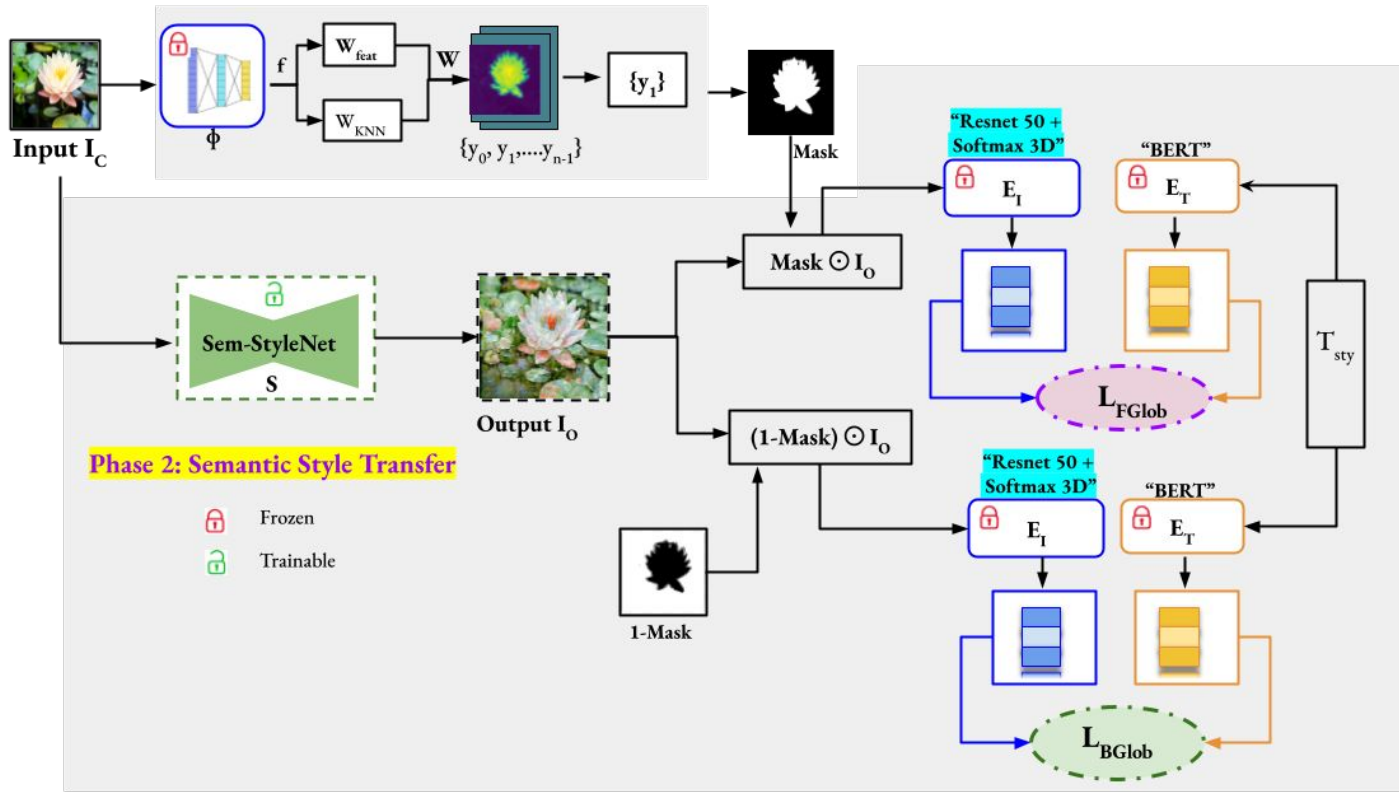
Connecting Text and Image through CLIP Space [5].

Background



Connecting Text and Image through CLIP Space [5].

Phase 1: Salient Object Detection



Notations

- I_C : Content Image
- T_{sty} : Style Text
- I_O : Stylized Output
- f : Deep patch features
- ϕ : Vision Transformer
- W_{KNN} : Color Matrix
- W_{feat} : Feature Matrix
- W : Semantic Affinity Matrix
- $\{y_0, y_1, \dots, y_{n-1}\}$: Eigen vectors
- S : Semantic StyleNet
- \odot : Hadamard Product
- E_I : CLIP Image Encoder
- E_T : CLIP Text Encoder
- L_{BGlob} : Global background loss
- L_{FGlob} : Global foreground loss