

Language-Driven Artistic Style Transfer



Tsu-Jui Fu¹

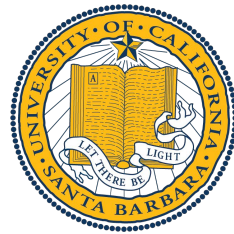


Xin Wang²



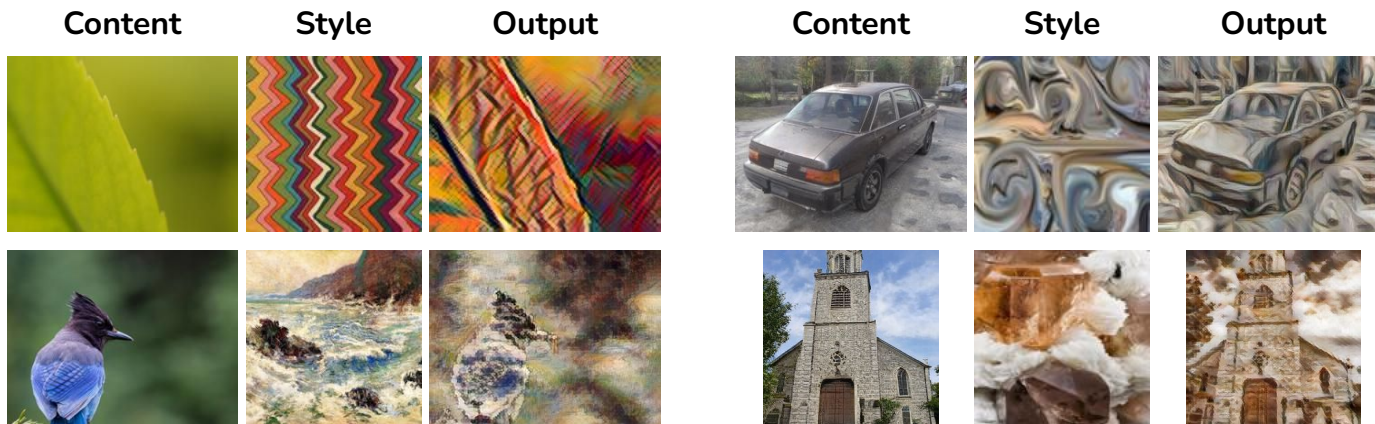
William Wang¹

¹UC Santa Barbara, ²UC Santa Cruz



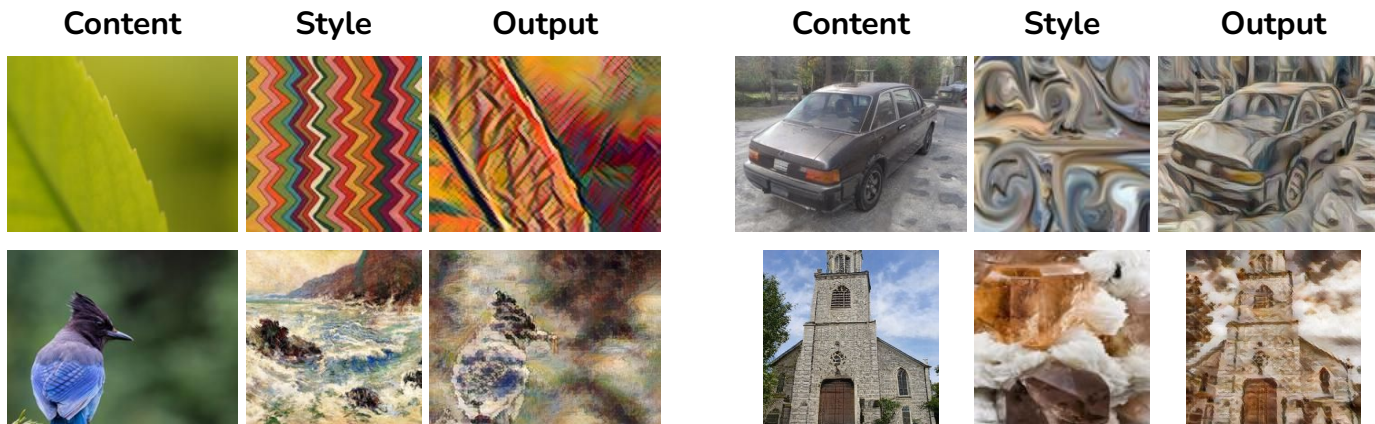
Artistic Style Transfer

- Render a photograph with an **arbitrary artwork style**
 - Preserve **content structures** yet present **style patterns**
- Content (\mathcal{C}) + Style (\mathcal{S}) \rightarrow Stylized Output (\mathcal{O})



Artistic Style Transfer

- Render a photograph with an **arbitrary artwork style**
 - Preserve **content structures** yet present **style patterns**
- Content (\mathcal{C}) + Style (\mathcal{S}) \rightarrow Stylized Output (\mathcal{O})



- Prepare collections of **style image** in advance
- Redraw new references **first** if there is no expected style

Language-Driven Artistic Style Transfer (LDAST)

- **Language** is the most natural way for humans to communicate
 - **Follow textual descriptions** to perform style transfer
 - Improve **accessibility** and **controllability**
- Content (\mathcal{C}) + Instruction (\mathcal{X}) \rightarrow Stylized Output (\mathcal{O})

Content



*out on a lovely day
with the water,
sketching, and painting*



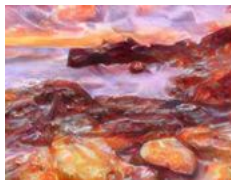
*i feel chaotic and
confused due to the
black and gray tones*



*peaceful green colors
and shading of the
branches, feel content*



*reflective, orange,
purple, and red bubble*



*salt deposits forming
around brown golden
frosted crystal*

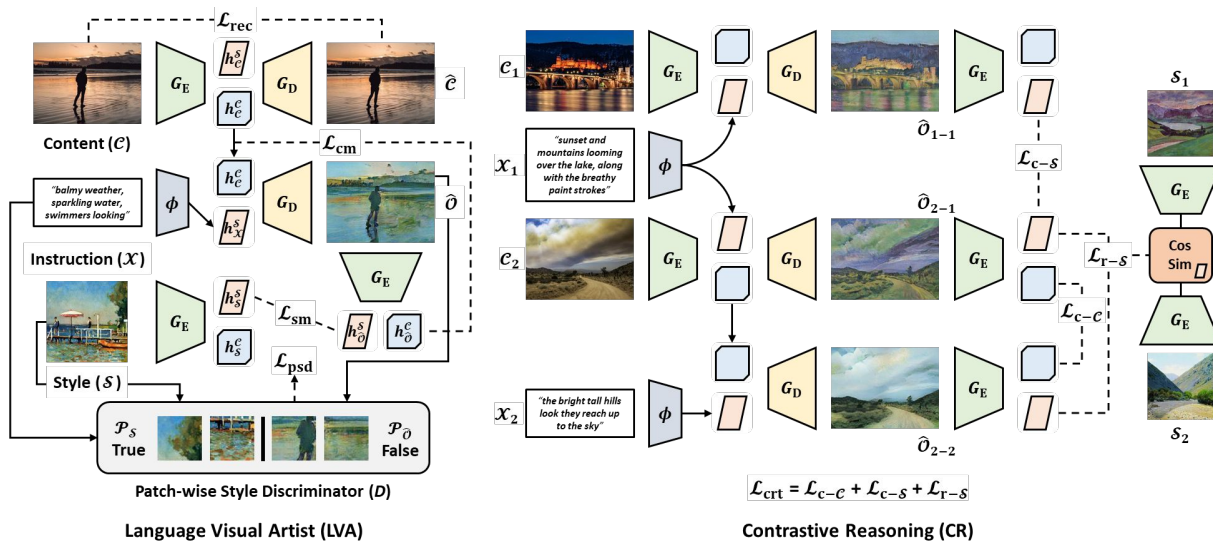


*large blackouts on
rough off white,
jute cotton surface*



Contrastive Language Visual Artist (CLVA)

- For **training**, there are content images (\mathcal{C}), style images (\mathcal{S}), and instructions (\mathcal{X})
- During **inference**, only \mathcal{C} and \mathcal{X} are provided
- Learn the **latent style patterns** from the instruction
- Further compare contrastive pairs of **relative \mathcal{C} and \mathcal{X}**

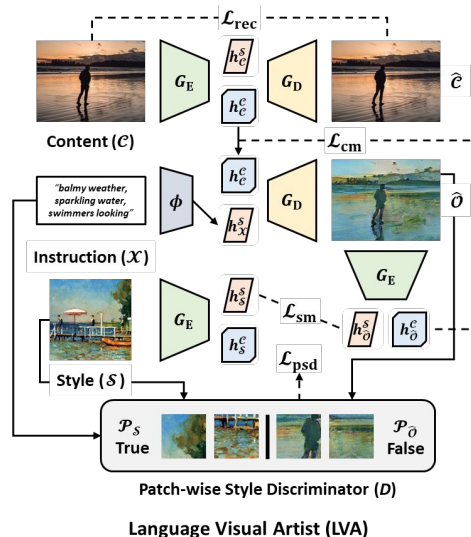


Language Visual Artist (LVA)

- Visual Encoder (\mathbf{G}_E), Text Encoder (Φ), and Visual Decoder (\mathbf{G}_D)
 - Extract content feature (h^C), style feature (h^S), and instruction feature (h^X)
 - Compose h^C and h^X / h^S to produce the stylized result
- **Structure Reconstruction (\mathcal{L}_{rec})**
 - Reproduce \mathcal{C} from the original content style
- **Patch-wise Style Discrimination (\mathcal{L}_{psd})**
 - D distinguishes the patch (\mathcal{P}) is from \mathcal{S} or \mathcal{O}
 - Optimize \mathbf{G}_E , Φ , and \mathbf{G}_D to fool D
- **Content Matching (\mathcal{L}_{cm}) and Style Matching (\mathcal{L}_{sm})**
 - Further enhance the alignment with the input

$$\mathcal{L}_{rec}, \mathcal{L}_{psd} = \|\hat{\mathcal{C}} - \mathcal{C}\|_2, \log(1 - D(\mathcal{P}_{\hat{\mathcal{O}}}, \mathcal{X})) + \log(D(\mathcal{P}_{\mathcal{S}}, \mathcal{X}))$$

$$\mathcal{L}_{cm}, \mathcal{L}_{sm} = \|h_{\hat{\mathcal{O}}}^C - h_{\mathcal{C}}^C\|_2, \|h_{\hat{\mathcal{O}}}^S - h_{\mathcal{S}}^S\|_2$$



Contrastive Reasoning (CR)

- Compare transferred results of **contrastive pairs** ($\{\mathcal{C}_1, \mathcal{X}_1, \mathcal{S}_1\}$ and $\{\mathcal{C}_2, \mathcal{X}_2, \mathcal{S}_2\}$)
 - Transfer to **various styles** while preserving the **same structure**
 - Apply **analogous style patterns** from **related style instructions**

- **Consistent Matching (\mathcal{L}_c)**

- Similar content structure from \mathcal{C}_2
- Similar style patterns from \mathcal{X}_1

- **Relative Matching (\mathcal{L}_r)**

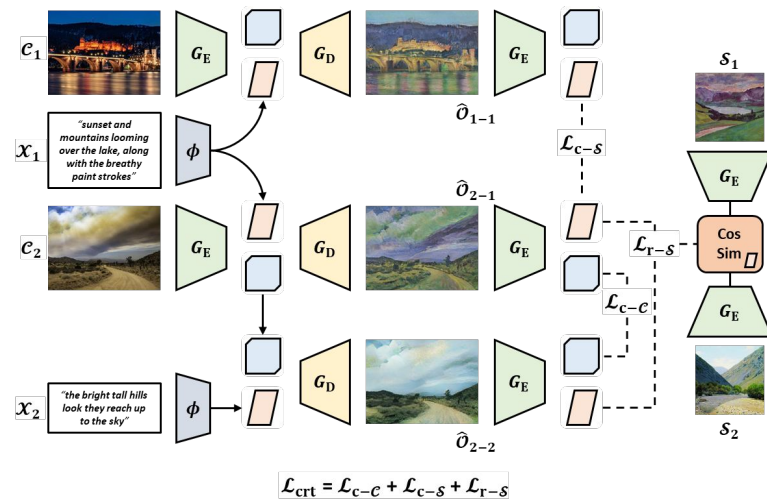
- Relative style patterns from \mathcal{X}_1 and \mathcal{X}_2

$$\mathcal{L}_{c-c} = \|h_{\hat{\mathcal{O}}_{c_1-x_1}}^c - h_{\hat{\mathcal{O}}_{c_1-x_2}}^c\|_2 + \|h_{\hat{\mathcal{O}}_{c_2-x_1}}^c - h_{\hat{\mathcal{O}}_{c_2-x_2}}^c\|_2$$

$$\mathcal{L}_{c-s} = \|h_{\hat{\mathcal{O}}_{c_1-x_1}}^s - h_{\hat{\mathcal{S}}_{2-1}}^s\|_2 + \|h_{\hat{\mathcal{O}}_{c_1-x_2}}^s - h_{\hat{\mathcal{O}}_{c_2-x_2}}^s\|_2$$

$$\mathcal{L}_{r-s} = (\|h_{\hat{\mathcal{O}}_{c_1-x_1}}^s - h_{\hat{\mathcal{O}}_{c_1-x_2}}^s\|_2 + \|h_{\hat{\mathcal{O}}_{c_2-x_1}}^s - h_{\hat{\mathcal{O}}_{c_2-x_2}}^s\|_2) \cdot r$$

$$\mathcal{L}_{ctr} = \mathcal{L}_{c-c} + \mathcal{L}_{c-s} + \mathcal{L}_{r-s}$$



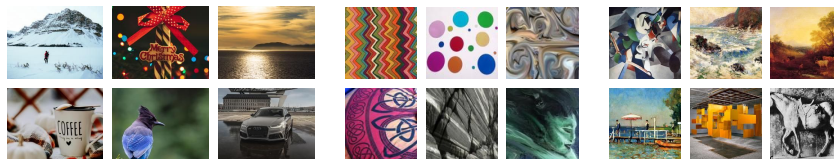
$$\mathcal{L}_{ctr} = \mathcal{L}_{c-c} + \mathcal{L}_{c-s} + \mathcal{L}_{r-s}$$

Contrastive Reasoning (CR)

Experimental Setup

- **Datasets**

- **Content:** Wallpaper
- **Style:** DTD² / ArtEmis



- **Evaluation Metrics (semi-GT from AdaAttN)**

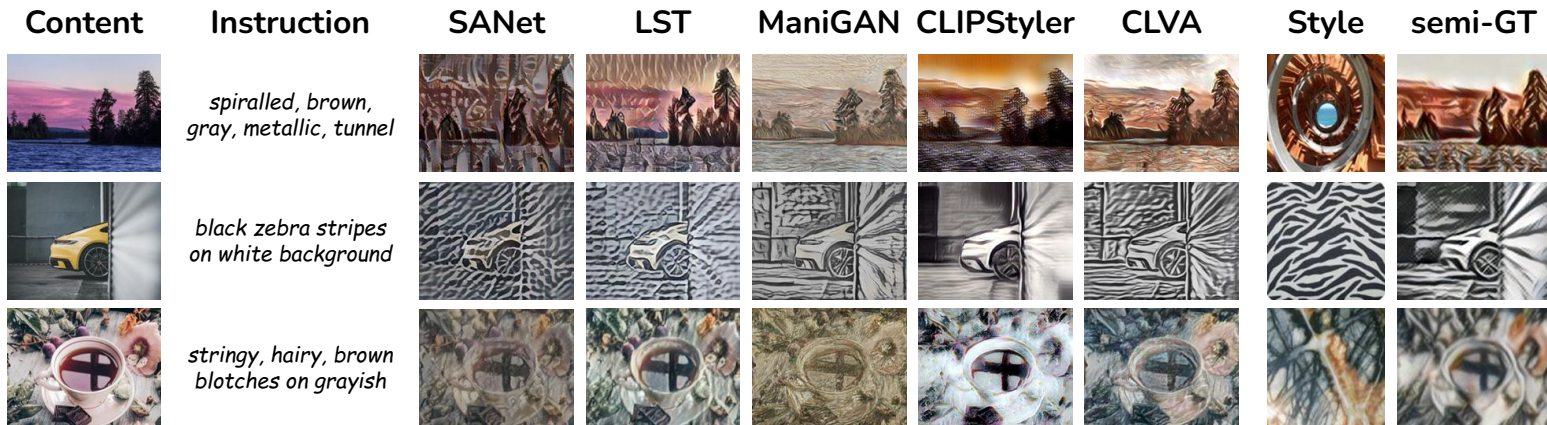
- **Percept** (↓): distance of **gram matrix** from visual features (vs. style image)
- **FAD** (↓): distance of **InceptionV3 features** (vs. semi-GT)
- **VLS** (↑): **relative visual-text similarity** from CLIP (vs. semi-GT | instruction)

- **Baselines**

- **Style Transfer:** SAnet / LST
- **Language-based Image Editing:** ManiGAN
- **CLIP-based Optimization:** StyleCLIP / NADA / CLIPStyler

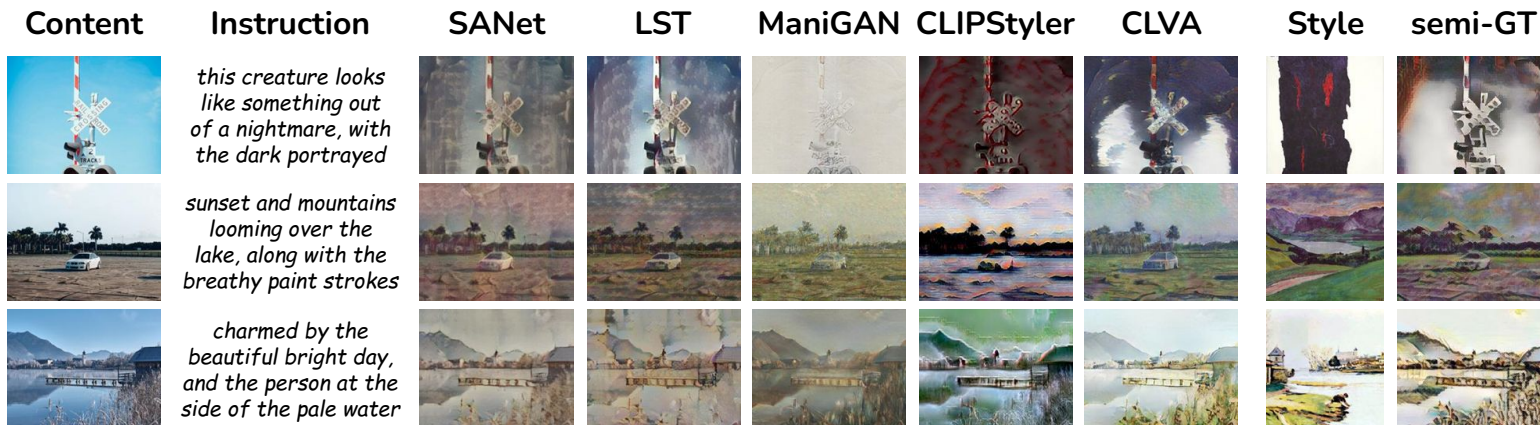
Instruction with Visual Attributes (DTD²)

Method	Automatic Metrics			Human Evaluation			
	Percept ↓	FAD ↓	VLS ↑	Content ↑	Instruction ↑	Style ↑	semi-GT ↑
SANet	<u>0.2129</u>	0.1627	23.57	2.701	2.477	2.738	2.630
LST	0.2129	<u>0.1533</u>	23.16	2.743	2.831	2.651	2.528
ManiGAN	0.2401	0.1663	23.25	2.757	2.562	2.937	2.922
CLIPStyler	0.2598	0.1818	24.62	<u>2.948</u>	<u>3.388</u>	<u>3.073</u>	<u>3.265</u>
CLVA	0.2033	0.1493	<u>24.00</u>	3.852	3.742	3.603	3.655



Instruction with Emotional Effects (ArtEmis)

Method	Automatic Metrics			Human Evaluation			
	Percept ↓	FAD ↓	VLS ↑	Content ↑	Instruction ↑	Style ↑	semi-GT ↑
SANet	0.0352	<u>0.1548</u>	19.30	<u>3.170</u>	2.978	2.980	2.890
LST	0.0386	0.1595	19.92	2.967	2.714	2.614	2.757
ManiGAN	0.0500	0.1554	19.69	2.729	2.583	2.879	<u>3.192</u>
CLIPStyler	0.0659	0.1759	21.04	2.777	<u>3.140</u>	<u>2.998</u>	2.952
CLVA	<u>0.0357</u>	0.1418	<u>20.11</u>	3.357	3.586	3.530	3.208



Specific Content Domain (Car & Church)

Method	Automatic Metrics			Human Evaluation			
	Percept ↓	FAD ↓	VLS ↑	Content ↑	Instruction ↑	Style ↑	semi-GT ↑
ManiGAN	<u>0.2329</u>	<u>0.1672</u>	23.44	2.861	2.894	2.978	2.893
StyleCLIP	0.2609	0.1812	21.55	3.459	2.845	2.930	2.829
NADA	0.2733	0.1876	23.38	2.542	2.798	2.846	2.932
CLIPStyler	0.2493	0.1826	24.16	2.986	<u>3.067</u>	<u>3.003</u>	<u>3.032</u>
CLVA	0.1957	0.1544	<u>23.68</u>	<u>3.153</u>	3.465	3.344	3.315



Ablation Study

- Reconstruction (\mathcal{L}_{rec}) + Patch-wise style (\mathcal{L}_{psd}) makes promising LDATAST
- Content matching (\mathcal{L}_{cm}) helps the **structure similarity**
- Style matching (\mathcal{L}_{sm}) aims at **analogous style patterns**
- Contrastive reasoning (\mathcal{L}_{ctr}) leads to a **comprehensive improvement**

$\mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{psd}}$	\mathcal{L}_{cm}	\mathcal{L}_{sm}	\mathcal{L}_{ctr}	Percept ↓	FAD ↓	VLS ↑
✓	✗	✗	✗	0.2290	0.1568	23.29
✓	✓	✗	✗	0.2304	0.1512	23.27
✓	✗	✓	✗	<u>0.2049</u>	0.1508	<u>23.69</u>
✓	✓	✓	✗	0.2100	<u>0.1499</u>	23.54
✓	✓	✓	✓	0.2033	0.1493	24.00

Why CLVA is better than CLIP-based?

- Investigate via **instruction-to-style retrieval**
 - CLIP cannot capture **detailed patterns** well

Method	DTD ²		ArtEmis	
	R@1	R@5	R@1	R@5
CLIP	13.9	30.7	9.8	20.7
CLVA	19.3	45.1	13.9	30.7

Method	Human Evaluation			
	Content ↑	Instruction ↑	Style ↑	semi-GT ↑
CLIPStyler (ft.)	1.208	1.347	1.292	1.333
CLVA	1.792	1.653	1.708	1.667

Instruction

all of the bright colors in the town makes it a happy place to live



lovely still life that looks like a tropical table setting



light green shiny embedded in a white rough and raised surface



Efficiency

- Evaluate on a **single TITAN X (12GB)** with content image size **256x192**
 - CLIP-based methods require **numerous iterations for optimization**
 - CLVA further takes advantage of **parallelization**

Method	Time (sec ↓)			GPU (MB ↓)		
	BS=1	32	50	BS=1	32	50
ManiGAN	0.079	0.533	1.148	3,312	6,572	8,129
StyleCLIP	32.38	*	*	4,149	*	*
NADA	63.49	*	*	6,413	*	*
CLIPStyler	99.98	*	*	5,429	*	*
CLVA	0.029	0.246	0.405	1,525	3,207	4,441






















* means this method can only run one input at a time

Linear Interpolation

- Consider two instructions \mathcal{X}_1 and \mathcal{X}_2
 - The **interpolated style feature** should be

$$h_p^S = (1 - \alpha)h_{\mathcal{X}_1}^S + \alpha h_{\mathcal{X}_2}^S$$

- Present a **smooth transformation** in between

Content	Instruction ₁	$\alpha=0.0$	$\alpha=0.2$	$\alpha=0.4$	$\alpha=0.6$	$\alpha=0.8$	$\alpha=1.0$	Instruction ₂
	<i>floating, colorful, white backdrop, circular round</i>							<i>transparent, white, brown, golden, rocky</i>
	<i>optical illusion with pen and ink drawing</i>							<i>the trees are very calming and warm</i>
	<i>transparent, white, brown, golden, rocky</i>							<i>the trees are very calming and warm</i>

Fine-grained Control

- Achieve fine-grained style control by **partial semantic editing**
 - The extracted patterns are **explicit** to reflect **each aspect of style semantic**



Super Resolution (2560x1440)

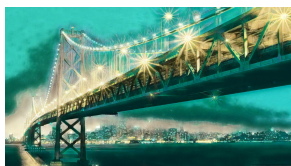
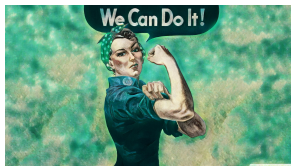
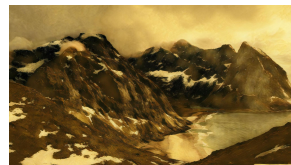
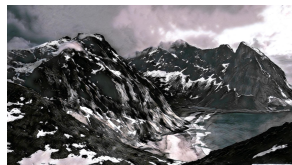
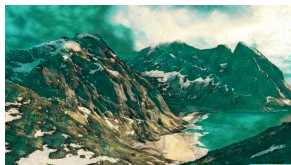
- Borrow from SAnet, which supports content images with **any resolutions**

*crystals like a
krypton diamond
in white and green*

*wrinkled leather,
gray shiny surface*

*painted, rubbed,
smeared yellow,
green, blue and red*

*warm painting
feels like the sun
is setting behind*



Conclusion

- Language-driven artistic style transfer (**LDAST**)
 - **Control artistic style transfer** via natural language
- Contrastive language visual artist (**CLVA**)
 - Learn to **extract explicit visual semantics** from style descriptions
 - Carry out instructions with **visual attributes / emotional effects**

*purple pink violet
medium polka dots*

*wrinkled, colorful,
soft fabric on
black background*

*sun is shining,
bouncing light,
summer scene*

*i feel chaotic and
confused due to the
black and gray tones*



Project



Code

