Visual Anagrams: Generating Multi-View Optical Illusions with Diffusion Models

CVPR 2024 (Oral)

Daniel Geng Inbum Park Andrew Owens University of Michigan

STRUCT Group Seminar Presenter: Lan Xicheng

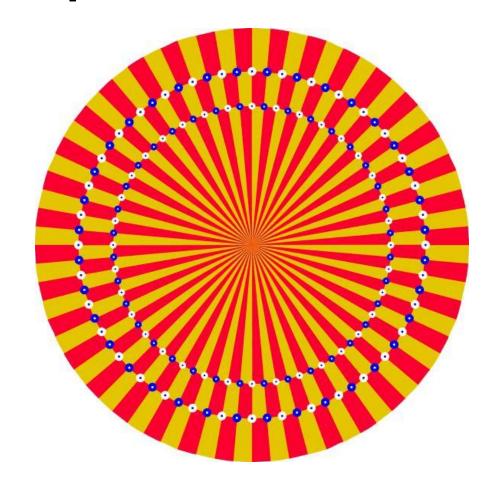
2024.9.14

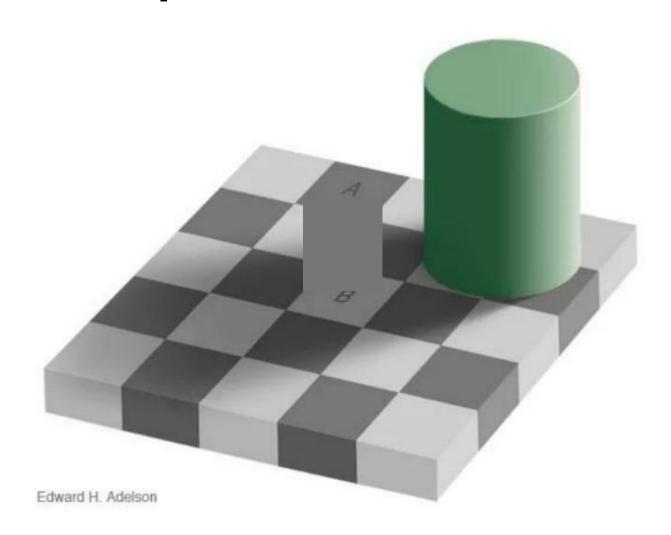
Outline

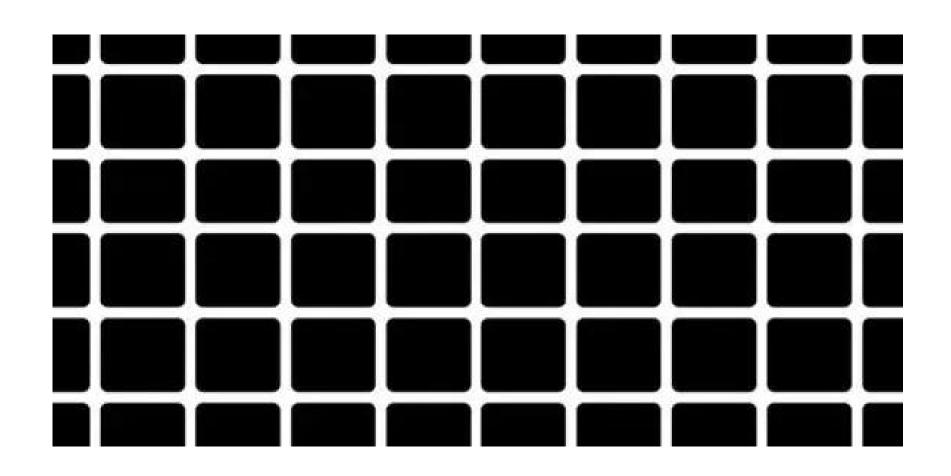
- Authors
- Background
- Methods
- Experiments
- Conclusion

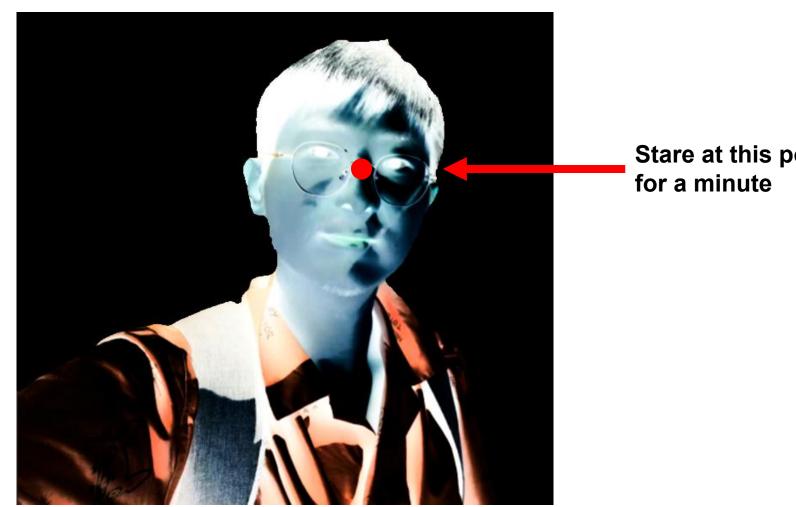
Outline

- Authors
- Background
- Methods
- Experiments
- Conclusion

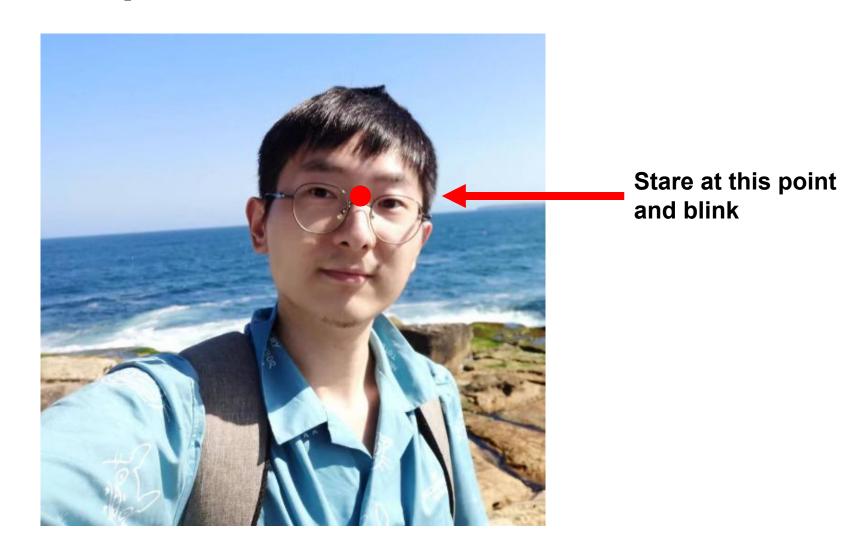






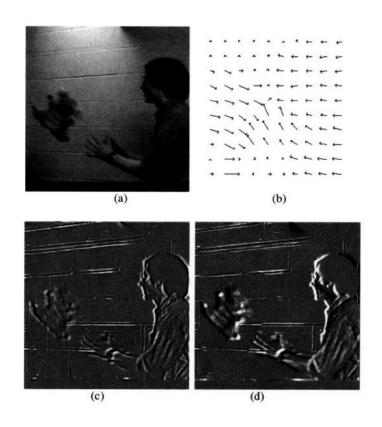


Stare at this point



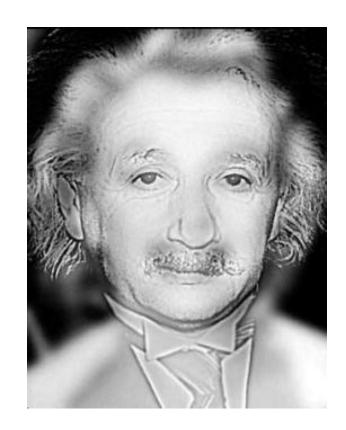
Computational Optical Illusions

- Motion Without Movement (SIGGRAPH 1991)
 Create the illusion of constant motion in a desired direction by locally applying a filter with continuously shifting phase
- Hybrid Images (TOG 2006)
- Camouflage Images (TOG 2010)
- Designing Perceptual Puzzles By Differentiating Probabilistic Programs (SIGGRAPH 2022)



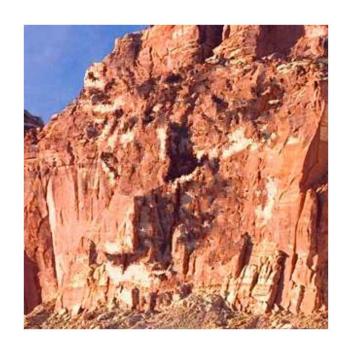
Computational Optical Illusions

- Motion Without Movement (SIGGRAPH 1991)
- Hybrid Images (TOG 2006)
 Change appearance depending on the distance they are viewed from
- Camouflage Images (TOG 2010)
- Designing Perceptual Puzzles By Differentiating Probabilistic Programs (SIGGRAPH 2022)



Computational Optical Illusions

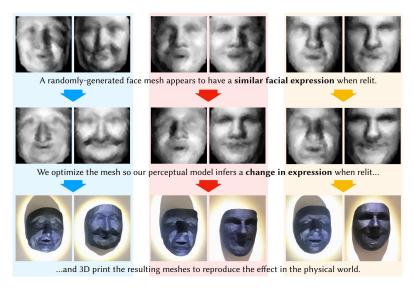
- Motion Without Movement (SIGGRAPH 1991)
- Hybrid Images (TOG 2006)
- Camouflage Images (TOG 2010)
 Camouflage objects in a scene through retexturing, with additional constraints on luminance as to preserve salient features of the object
- Designing Perceptual Puzzles By Differentiating Probabilistic Programs (SIGGRAPH 2022)



Computational Optical Illusions

- Motion Without Movement (SIGGRAPH 1991)
- Hybrid Images (TOG 2006)
- Camouflage Images (TOG 2010)
- Designing Perceptual Puzzles By Differentiating Probabilistic Programs (SIGGRAPH 2022)

Color-constancy, size constancy, and face perception illusions by differentiating through a Bayesian model of human vision



Illusions with Diffusion Models

- QR Codes As Created Images
 Global structure subtly matches a given template image (I2I)
- Extensions



Illusions with Diffusion Models

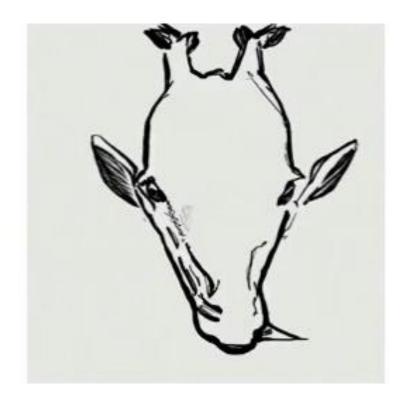
- QR Codes As Created Images
- Extensions



Illusions with Diffusion Models

- QR Codes As Created Images
- Extensions





a drawing of a giraffe



an oil painting of a snowy mountain village



an oil painting of a skull



a pop art of marilyn monroe

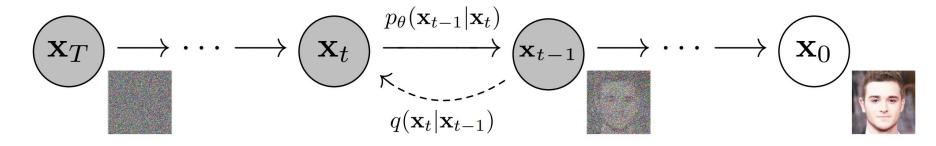


a lithograph of a mountain landscape



a watercolor of a kitten

Diffusion Models



Algorithm 1 Training

- 1: repeat
- 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on

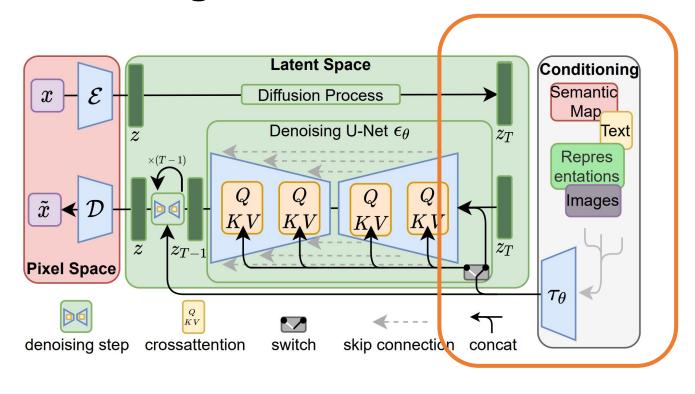
$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} (\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$

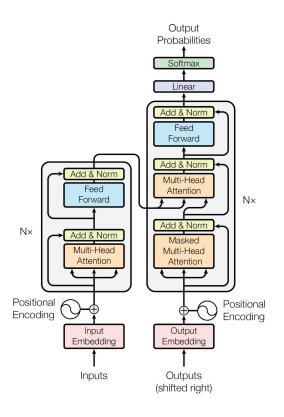
6: **until** converged

Algorithm 2 Sampling

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for** t = T, ..., 1 **do**
- 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if t > 1, else $\mathbf{z} = \mathbf{0}$
- 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
- 5: end for
- 6: return x_0

Text-to-Image





Compositional Generation

- Composable Diffusion Models (ECCV 2022)
- Reduce, Reuse, Recycle (PMLR 2023)

Composing Language Descriptions (Composed Stable Diffusion)



"A photo of cherry blossom trees" AND "Sun dog" AND "Green grass"



"A church" AND "Lightning in the background" AND "A beautiful pink sky"



"A stone castle and trees," AND "Black and white"



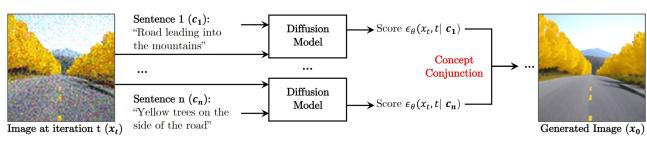
"A stone castle surrounded by lakes surrounded by lakes and trees," AND (NOT "Black and white")



"A mystical tree" AND "A dark magical pond" AND "Dark"



"A mystical tree" AND "A dark magical pond" AND (NOT "Dark")



Compositional Generation

- Composable Diffusion Models (ECCV 2022)
- Reduce, Reuse, Recycle (PMLR 2023)

$$p_{\theta}(\boldsymbol{x}) \propto e^{-E_{\theta}(\boldsymbol{x})}$$
$$\boldsymbol{x}_{t} = \boldsymbol{x}_{t-1} - \frac{\lambda}{2} \nabla_{\boldsymbol{x}} E_{\theta}(\boldsymbol{x}_{t-1}) + \mathcal{N}(0, \sigma_{t}^{2} I)$$

Energy-Based Models (EBMs)

$$p_{\text{compose}}(\boldsymbol{x}) \propto p_{\theta}^{1}(\boldsymbol{x}) \cdots p_{\theta}^{n}(\boldsymbol{x}) \propto e^{-\sum_{i=1}^{n} E_{\theta}^{i}(\boldsymbol{x})}$$
$$\boldsymbol{x}_{t} = \boldsymbol{x}_{t-1} - \frac{\lambda}{2} \nabla_{\boldsymbol{x}} \left(\sum_{i=1}^{n} E_{\theta}^{i}(\boldsymbol{x}_{t-1}) \right) + \mathcal{N}(0, \sigma_{t}^{2}I)$$

Composing EBMs

$$p_{\text{compose}}(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t) = \mathcal{N}\left(\boldsymbol{x}_t - \sum_{i=1}^n \epsilon_{\theta}^i(\boldsymbol{x}_t, t), \sigma_t^2 I\right)$$

Composing Diffusion Models

: Image $E_{ heta}(oldsymbol{x})$: Energy Function,

 \boldsymbol{x}

Learnable

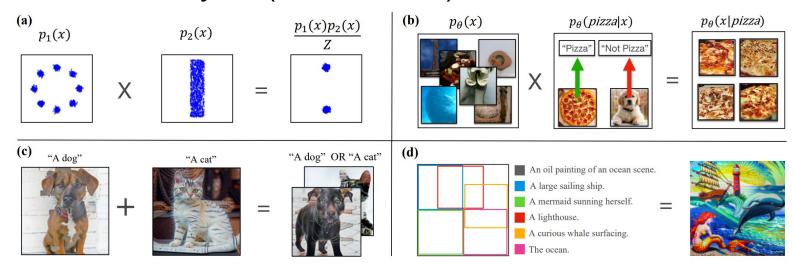
 $\epsilon_{ heta}^{i}$: Diffusion Model

 $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) \coloneqq \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$

24

Compositional Generation

- Composable Diffusion Models (ECCV 2022)
- Reduce, Reuse, Recycle (PMLR 2023)

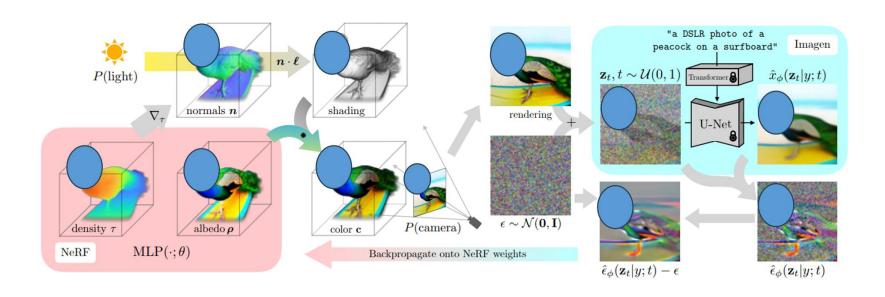


More theoretical, More in-depth, Not just relying on diffusion

Illusions with Diffusion Models

Baseline: Score Distillation Sampling (SDS)

Create images align with different prompts from different views



Illusions with Diffusion Models

Baseline: Matthew Tancik. Illusion diffusion. https://github.com/tancik/Illusion-Diffusion

- Combining noise predictions from different views during denoising
- Based on latent diffusion model (LDM)
- Just for rotation illusions



Outline

- Authors
- Background
- Methods
- Experiments
- Conclusion

Text-conditioned Diffusion Models

Classifier-free Guidance (CFG)

$$\epsilon_t^{\text{CFG}} = \epsilon_{\theta}(\mathbf{x}_t, t, \varnothing) + \gamma(\epsilon_{\theta}(\mathbf{x}_t, t, y) - \epsilon_{\theta}(\mathbf{x}_t, t, \varnothing))$$

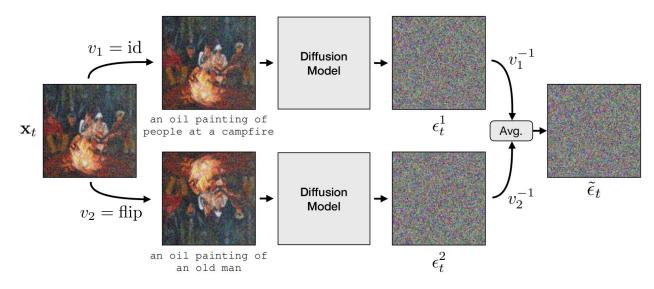
y: Conditioning Prompts \varnothing : Embedding Of The Empty String γ : Strength Parameter

Negative prompting: the empty text prompt embedding, \emptyset , is replaced by a text prompt that discourage the model from generating.

Parallel Denoising

Simultaneously Denoise Multiple Views Of An Image

$$\tilde{\epsilon}_t = \frac{1}{N} \sum_i v_i^{-1} \left(\epsilon_{\theta}(v_i(\mathbf{x}_t), y_i, t) \right)$$



Parallel Denoising

Simultaneously Denoise Multiple Views Of An Image

$$ilde{\epsilon}_t = rac{1}{N} \sum_i v_i^{-1} \left(\epsilon_{ heta}(v_i(\mathbf{x}_t), y_i, t)
ight)$$
Replace as $\epsilon_t^{ ext{CFG}}$

N: Number Of Prompt Set y_i : Prompt i

 $v_i(\cdot)$: View Function i $\epsilon_{ heta}(\cdot)$: Diffusion Model

$$p_{ ext{compose}}(m{x}_{t-1}|m{x}_t) = \mathcal{N}\left(m{x}_t - \sum_{i=1}^n \epsilon_{ heta}^i(m{x}_t,t), \sigma_t^2 I
ight)$$
 Composing Diffusion Models

Conditions On Views

- Invertibility
- Linearity
- Statistical Consistency

Conditions On Views

- Invertibility
- Linearity

$$\mathbf{x}_t = w_t^{\text{signal}} \underbrace{\mathbf{x}_0}_{\text{signal}} + w_t^{\text{noise}} \underbrace{\epsilon}_{\text{noise}}$$

 $v_i(\mathbf{x}_t)$ also need to be a linear combination of pure signal and pure noise with the same weighting

$$v_i(\mathbf{x}_t) = \mathbf{A}_i(w_t^{\text{signal}}\mathbf{x}_0 + w_t^{\text{noise}}\epsilon)$$

$$= w_t^{\text{signal}} \underbrace{\mathbf{A}_i\mathbf{x}_0}_{\text{new signal}} + w_t^{\text{noise}} \underbrace{\mathbf{A}_i\epsilon}_{\text{new noise}}$$

Statistical Consistency

Conditions On Views

- Invertibility
- Linearity
- Statistical Consistency
 - -Diffusion model $\epsilon \sim \mathcal{N}(0,I)$
 - -Transformed noise $\mathbf{A}_i\epsilon$ must be likewise drawn from $\mathcal{N}(0,I)$
 - - \mathbf{A}_i is an orthogonal matrix

Views Considered

Standard Image Manipulations

Rotation, reflection and skewing



an oil painting of an old man



an oil painting of a skull

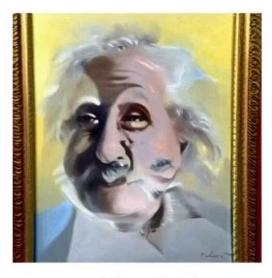
Views Considered

General Permutations

Jigsaw puzzles and inner rotations



an oil painting of a deer



an oil painting of albert einstein

Views Considered

Color Inversion



a lithograph of a teddy bear



a photo of a woman

Views Considered

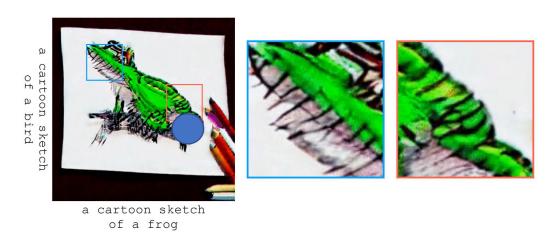
Arbitrary Orthogonal Transformations



Design Decisions

Pixel Diffusion Model

Latent diffusion models lead to artifacts under rotations or flips, where the location of latents change, but the content and orientation of these blocks do not



- Combining Noise Estimates
- Negative Prompting

Design Decisions

- Pixel Diffusion Model
- Combining Noise Estimates

Alternating through them by timestep

$$\tilde{\epsilon}_t = v_{t \bmod N}^{-1} \left(\epsilon_{\theta}(v_{t \bmod N}(\mathbf{x}_t), t, y) \right)$$

Negative Prompting

Design Decisions

- Pixel Diffusion Model
- Combining Noise Estimates
- Negative Prompting

Use one view's prompt as a negative for the other view, and vice versa This encourages the model to hide the other view's prompt for a given view

Outline

- Authors
- Background
- Methods
- Experiments
- Conclusion

Metrics

- CLIP to measure how well views align with the desired prompts
- $\mathbf{S} \in \mathbb{R}^{N \times N}$ defined as $\mathbf{S}_{ij} = \phi_{\mathrm{img}}(v_i(\mathbf{x}))^T \phi_{\mathrm{text}}(p_j)$ ϕ_{img} and ϕ_{text} are the CLIP visual and textual encoders respectively
- Alignment score: $\mathcal{A} = \min \operatorname{diag}(\mathbf{S})$, measures the worst alignment
- Concealment score: $\mathcal{C}=\frac{1}{N}\operatorname{tr}(\operatorname{softmax}(S/\tau))$, measures how well CLIP can classify a view

Datasets: Prompt Pairs For 2-view Illusions

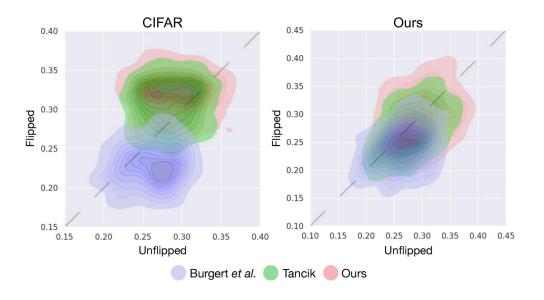
- CIFAR: 10 classes from CIFAR-10, for a total of 45 prompt pairs
- Ours: compile by hand, consists of 50 prompt pairs

Baselines

- Burgert et al. : Score Distillation Sampling
- Tancik: An earlier version of our method

Quantitative Results

Prompt Pair	Method	$\mathcal{A}\uparrow$	$\mathcal{A}_{0.9}\uparrow$	$\mathcal{A}_{0.95}\uparrow$	$\mathcal{C}\uparrow$	$\mathcal{C}_{0.9}\uparrow$	$\mathcal{C}_{0.95} \uparrow$
CIFAR	Burgert <i>et al</i> . [2] Tancik [42] Ours		0.310	0.260 0.316 0.327	0.595		
Ours	Burgert <i>et al</i> . [2] Tancik [42] Ours	0.256	0.270 0.294 0.315	0.283 0.309 0.326	0.545	0.621	0.655



Qualitative Results

a painting of a truck

a painting
 of a deer

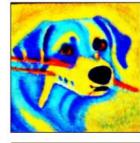
a painting a painting of a dog of an airplane

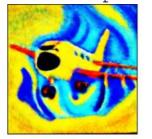
an ink drawing an ink drawing of a house of a castle

Burgert et al.

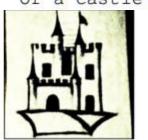












Tancik



























Ablations

Ablation	$\mathcal{A}\uparrow$	$\mathcal{A}_{0.9}\uparrow$	$\mathcal{A}_{0.95}\uparrow$	$\mathcal{C}\uparrow$	$\mathcal{C}_{0.9}\uparrow$	$\mathcal{C}_{0.95} \uparrow$
Negative Prompting No Negative Prompting	0.24 0.255	0.27 0.285	0.276 0.295	0.576 0.567	0.659 0.643	0.683 0.679
Alternating Reduction Mean Reduction	0.252 0.255	0.286 0.285	0.292 0.295	0.560 0.567	0.639 0.643	0.664 0.679
$\begin{array}{c} \gamma = 3.0 \\ \gamma = 7.0 \end{array}$	0.239 0.255	0.271 0.285	0.285 0.295	0.537 0.567	0.610 0.643	0.629 0.679
$\gamma = 10.0$	0.259	0.290	0.297	0.576	0.664	0.702

Outline

- Authors
- Background
- Methods
- Experiments
- Conclusion

Conclusions

- Present a method to produce compelling and diverse optical illusions
- Prove the method works for a broad set of transformations
- Qualitatively show the method can generate a wide array of optical illusions

Does not consistently produce perfect illusions

Failures



Failure: Independent Synthesis

View: Vertical Flip

an oil painting of a bowl of fruit

an oil painting of a monkey



Failure: Noise Shift

View: White Balancing

a photo of a black and blue dress a photo of a white and gold dress



Failure: Correlated Noise

View: Rotate Inner Circle 45 Degrees (Bilinear)

a sketch of an elephant

a sketch of a mouse

Thank you for Listening!