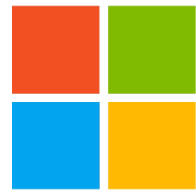




西安交通大学  
XI'AN JIAOTONG UNIVERSITY



Microsoft

# Diffusion Model with Cross Attention as an Inductive Bias for Disentanglement

Speaker: Tao Yang



西安交通大学  
XI'AN JIAOTONG UNIVERSITY



# CONTENTS

1. Background

2. Method

3. Experiment

## (4) Disentangled Models & Visual Concept Learning

Objective: Empower the existing AI models the ability of human-like concept induction, and develop a unified foundation model to achieve artificial general intelligence.

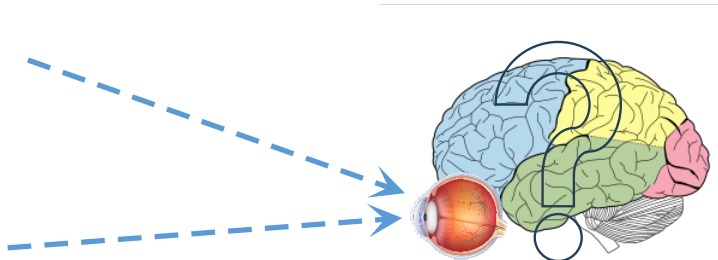


An AI must fundamentally understand the world around us, and we argue that this can only be achieved if it can learn to identify and disentangle the underlying explanatory factors hidden in the observed milieu of low-level sensory data. —Geoffrey Hinton “Distributed Representation”  
Concepts can be represented by distributed patterns of activity in networks of neuron-like units. —Yoshua Bengio “Representation Learning: A

# 1. Background

## (5) Challenges

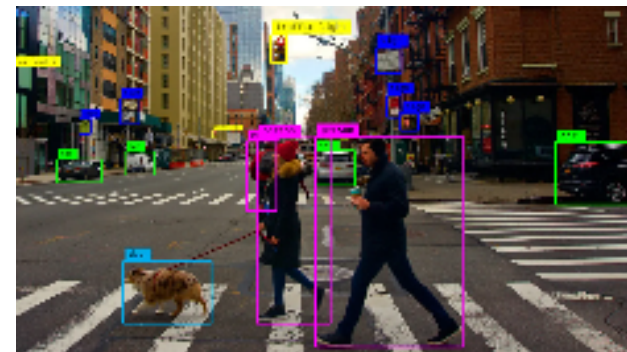
### □ Unknown Mechanism of Induction



Color: black, brown, ...  
Shape: square, round, ...  
Material: wooden, ...  
Azimuth: left, right, ...

### □ Complexity of Image Data

- Occlusion
- Complex Concept
- Complex Scene

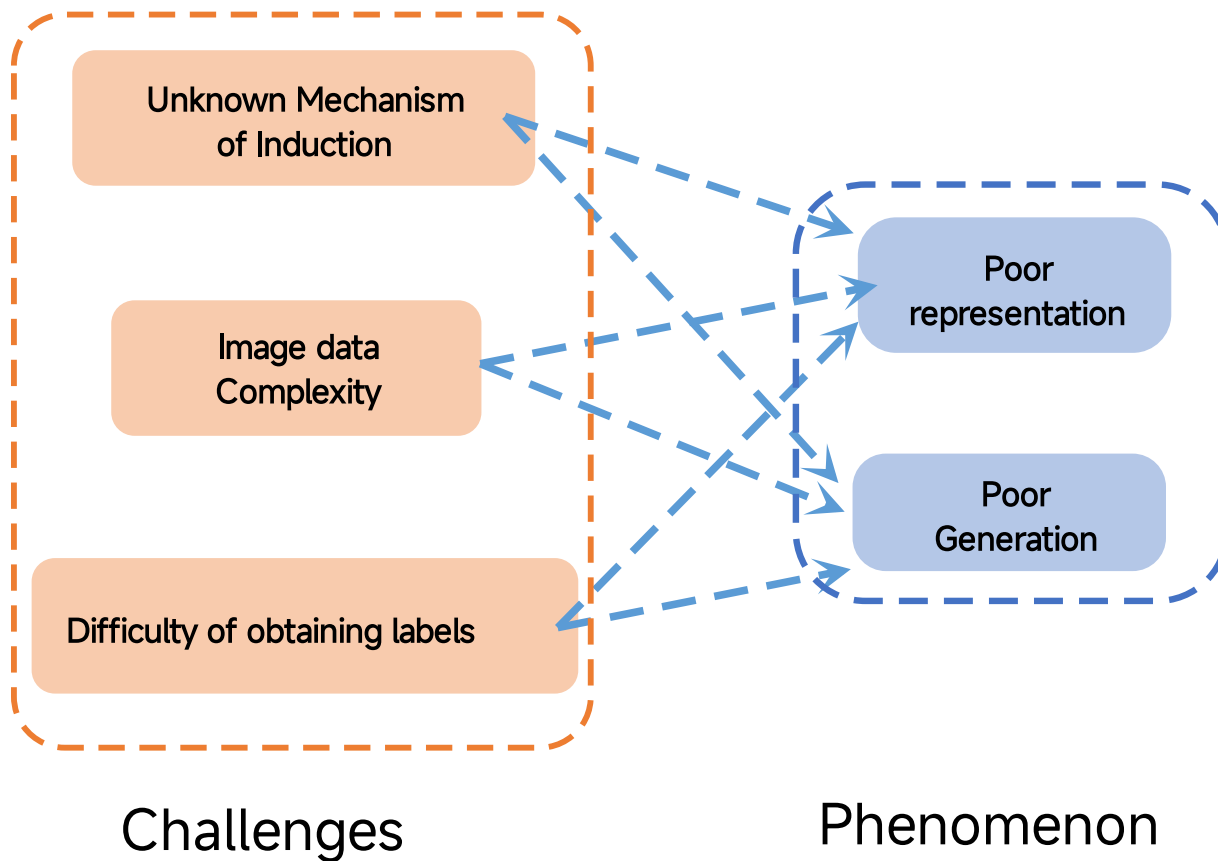


### □ Difficulty of obtaining labels: there are many concepts and many objects

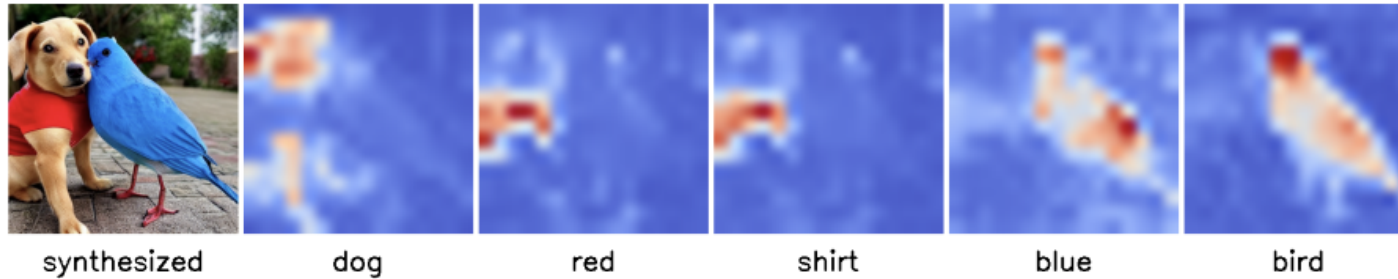


# 1. Background

## (5) Challenges



## 2. Visual Concept Learning by Disentanglement



### □ Motivation

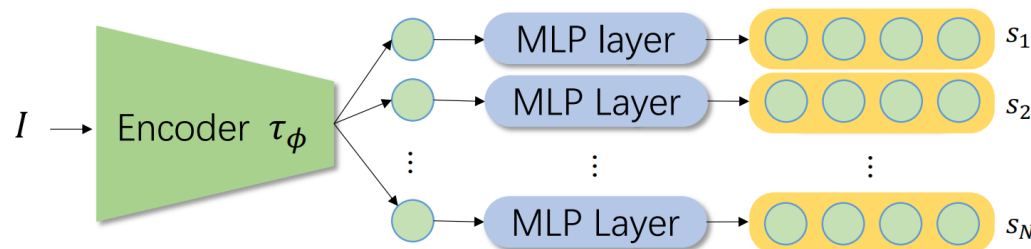
- VAEs and GANs based methods primarily rely on probability-based regularizations applied to the latent space.
- In text-to-image generation, a conditional diffusion model integrates the “disentangled” text tokens by cross attention, demonstrating the ability to generate semantically aligned images.
- Locatello et al. (2019) demonstrate that relying solely on regularizations is insufficient

## 2. Visual Concept Learning by Disentanglement

### □ Network Design

- The image encoder  $\tau_\phi$  aims to provide a set of concept tokens  $S = \{s_1, \dots, s_N\}$ , which act similarly to the word embeddings in the prompts for text-to-image generation in stable diffusion.

- We treat each dimension of the encoded feature vector as a disentangled factor and map each factor to a vector (i.e., concept token) by non-shared MLP layers



## 2. Visual Concept Learning by Disentanglement

### □ Network Design

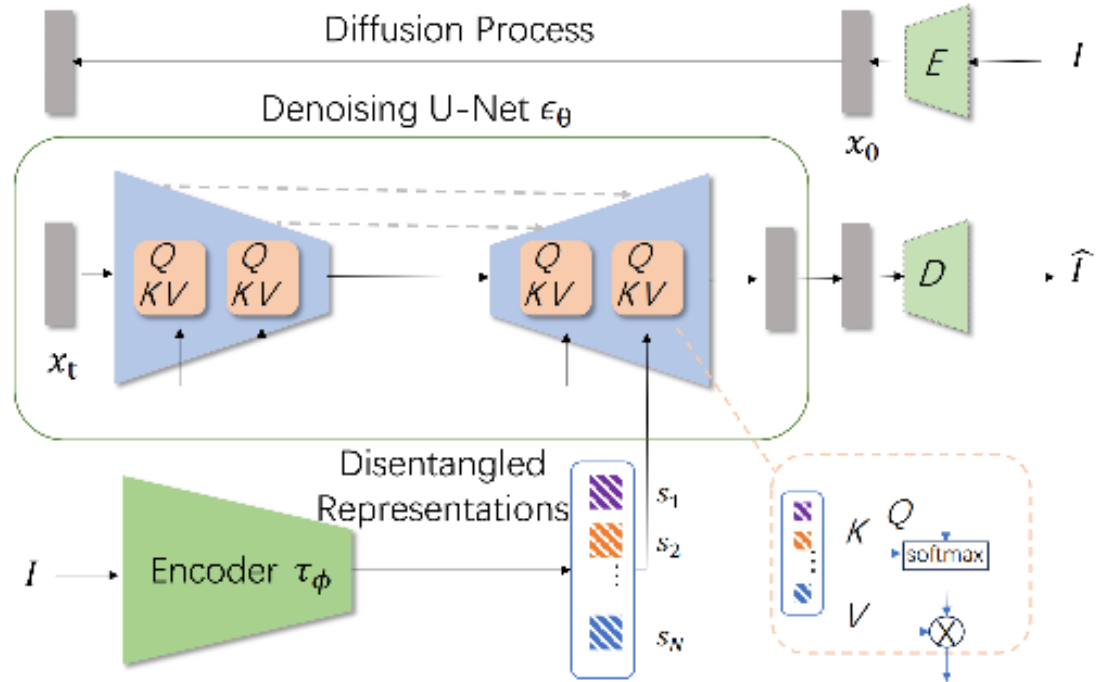
To condition the concept tokens during image generation, cross-attention is used to map these tokens into the representations of the U-Net in the diffusion model. This is

accomplished using the cross attention mechanism

$$Attention(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right) \cdot V$$

where each spatial feature in the intermediate feature map in diffusion serves as a query, the concept tokens act as keys

and values.



(a) Framework of our EncDiff.

## 2. Visual Concept Learning by Disentanglement

### □ The Information Bottleneck

The diffusion model optimizes a network (e.g., U-Net)  $\epsilon_\theta$  to predict the noise from the noisy input  $x_t$  and the condition  $\mathcal{S}$  (concept tokens), with the loss function defined as

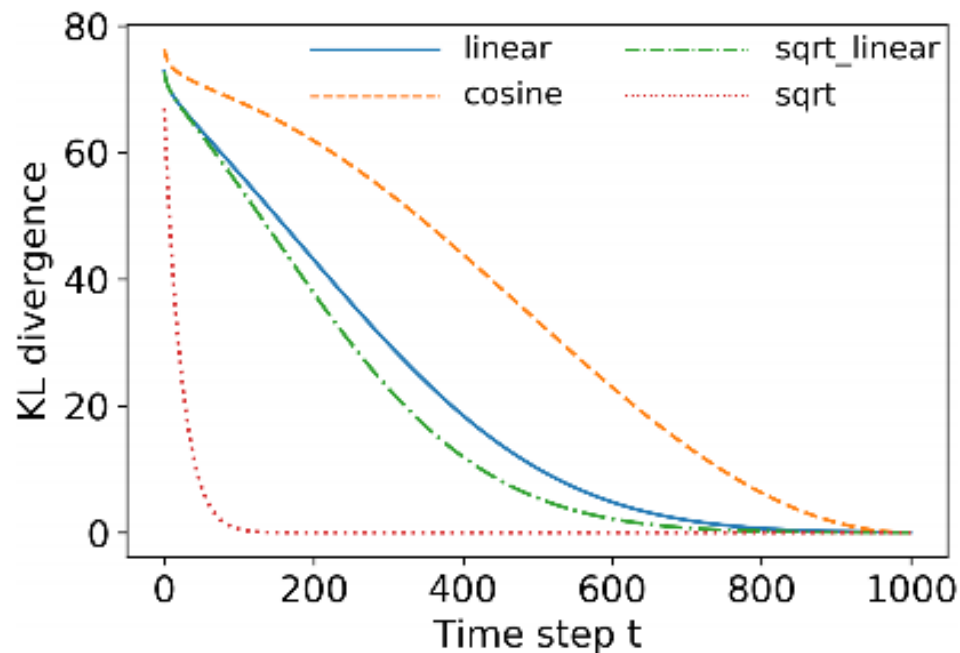
$$\mathcal{L}_r = \mathbb{E}_{x_0, \epsilon, t} \|\epsilon_\theta(x_t, t, \mathcal{S}) - \epsilon\|.$$

We can interpret the loss function as a reconstruction of latent input  $x_0$ . The loss function  $\mathcal{L}_r$  can be rewritten as:

$$\mathcal{L}_r = \sum_{t \geq 1} |C_t - D(p_\theta(x_{t-1}|x_t, \mathcal{S}) || q(x_{t-1}))|,$$

The loss function of AnnealVAE

$$\begin{aligned} \mathcal{L}(\varphi, \phi) = & -E_{q_\varphi(z|x)}[\log p_\phi(x|z)] \\ & + \gamma |C - D_{KL}(q_\varphi(z|x) || p(z))|, \end{aligned}$$



(b) KL divergence curves.



## 2. Visual Concept Learning by Disentanglement

### Quantitative Analysis

Method	Cars3D		Shapes3D		MPI3D	
	FactorVAE score $\uparrow$	DCI $\uparrow$	FactorVAE score $\uparrow$	DCI $\uparrow$	FactorVAE score $\uparrow$	DCI $\uparrow$
<i>VAE-based:</i>						
FactorVAE [17]	0.906 $\pm$ 0.052	0.161 $\pm$ 0.019	0.840 $\pm$ 0.066	0.611 $\pm$ 0.082	0.152 $\pm$ 0.025	0.240 $\pm$ 0.051
$\beta$ -TCVAE [4]	0.855 $\pm$ 0.082	0.140 $\pm$ 0.019	0.873 $\pm$ 0.074	0.613 $\pm$ 0.114	0.179 $\pm$ 0.017	0.237 $\pm$ 0.056
<i>GAN-based:</i>						
InfoGAN-CR [21]	0.411 $\pm$ 0.013	0.020 $\pm$ 0.011	0.587 $\pm$ 0.058	0.478 $\pm$ 0.055	0.439 $\pm$ 0.061	0.241 $\pm$ 0.075
<i>Pre-trained GAN-based:</i>						
LD [29]	0.852 $\pm$ 0.039	0.216 $\pm$ 0.072	0.805 $\pm$ 0.064	0.380 $\pm$ 0.062	0.391 $\pm$ 0.039	0.196 $\pm$ 0.038
GS [11]	0.932 $\pm$ 0.018	0.209 $\pm$ 0.031	0.788 $\pm$ 0.091	0.284 $\pm$ 0.034	0.465 $\pm$ 0.036	0.229 $\pm$ 0.042
DisCo [26]	0.855 $\pm$ 0.074	0.271 $\pm$ 0.037	0.877 $\pm$ 0.031	0.708 $\pm$ 0.048	0.371 $\pm$ 0.030	0.292 $\pm$ 0.024
<i>Diffusion-based:</i>						
DisDiff [37]	<b>0.976</b> $\pm$ 0.018	0.232 $\pm$ 0.019	0.902 $\pm$ 0.043	0.723 $\pm$ 0.013	0.617 $\pm$ 0.070	0.337 $\pm$ 0.057
EncDiff (Ours)	0.773 $\pm$ 0.060	<b>0.279</b> $\pm$ 0.022	<b>0.999</b> $\pm$ 0.000	<b>0.969</b> $\pm$ 0.030	<b>0.872</b> $\pm$ 0.049	<b>0.685</b> $\pm$ 0.044

## 2. Visual Concept Learning by Disentanglement

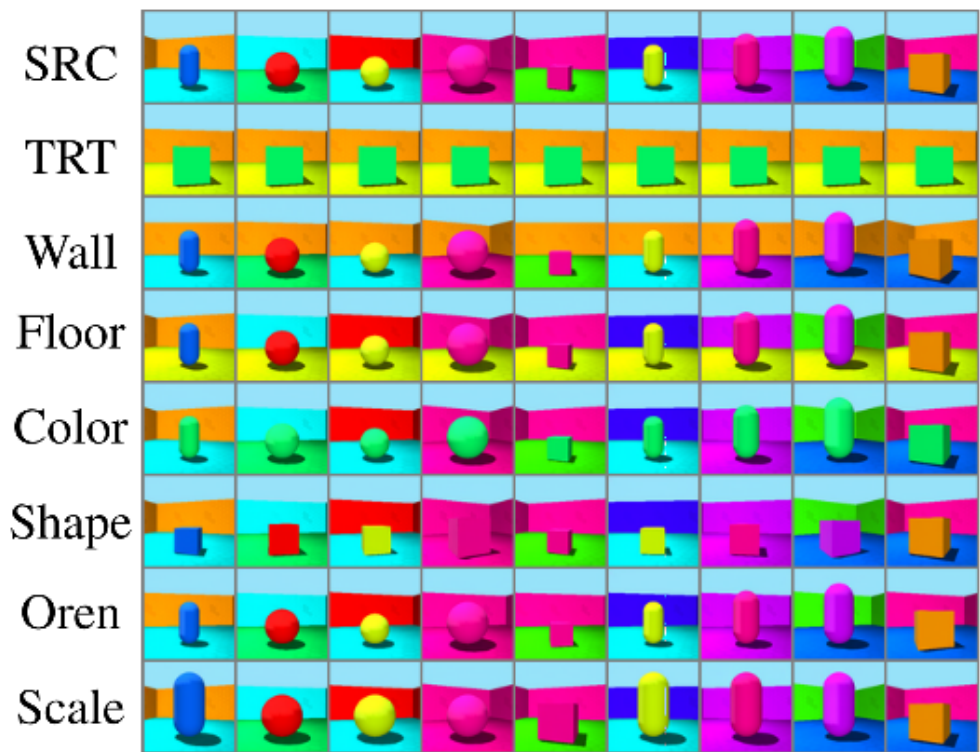


### Quantitative Results on CelebA

Model	TAD $\uparrow$	FID $\downarrow$
$\beta$ -VAE [14]	$0.088 \pm 0.043$	$99.8 \pm 2.4$
InfoVAE [40]	$0.000 \pm 0.000$	$77.8 \pm 1.6$
Diff-AE [24]	$0.155 \pm 0.010$	$22.7 \pm 2.1$
InfoDiffusion [31]	$0.299 \pm 0.006$	$23.6 \pm 1.3$
DisDiff [37]	$0.305 \pm 0.010$	$18.2 \pm 2.1$
EncDiff	<b><math>0.638 \pm 0.008</math></b>	<b><math>14.8 \pm 2.3</math></b>

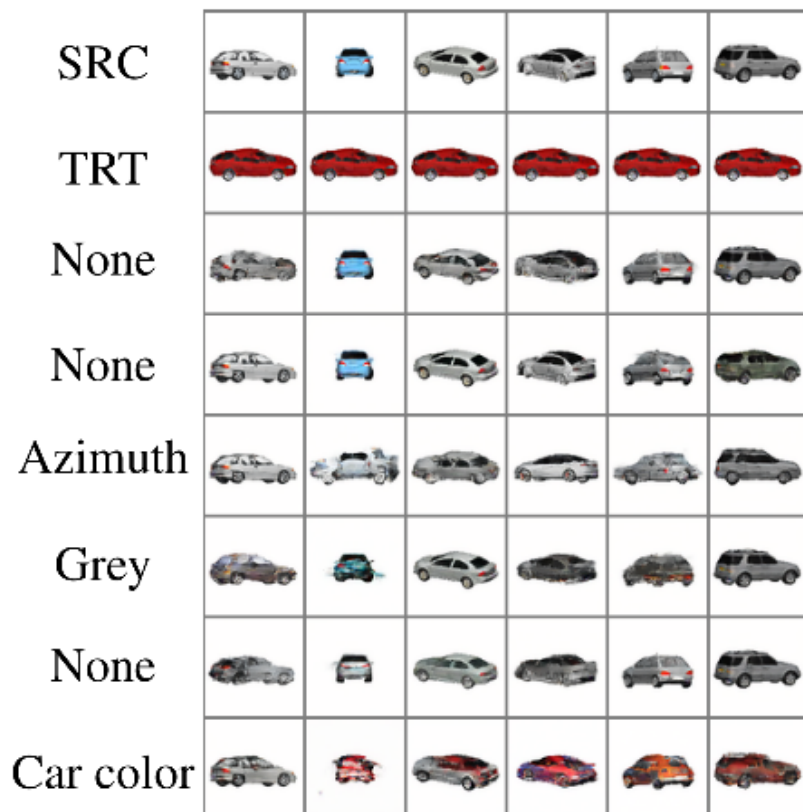
## 2. Visual Concept Learning by Disentanglement

### □ Qualitative Results

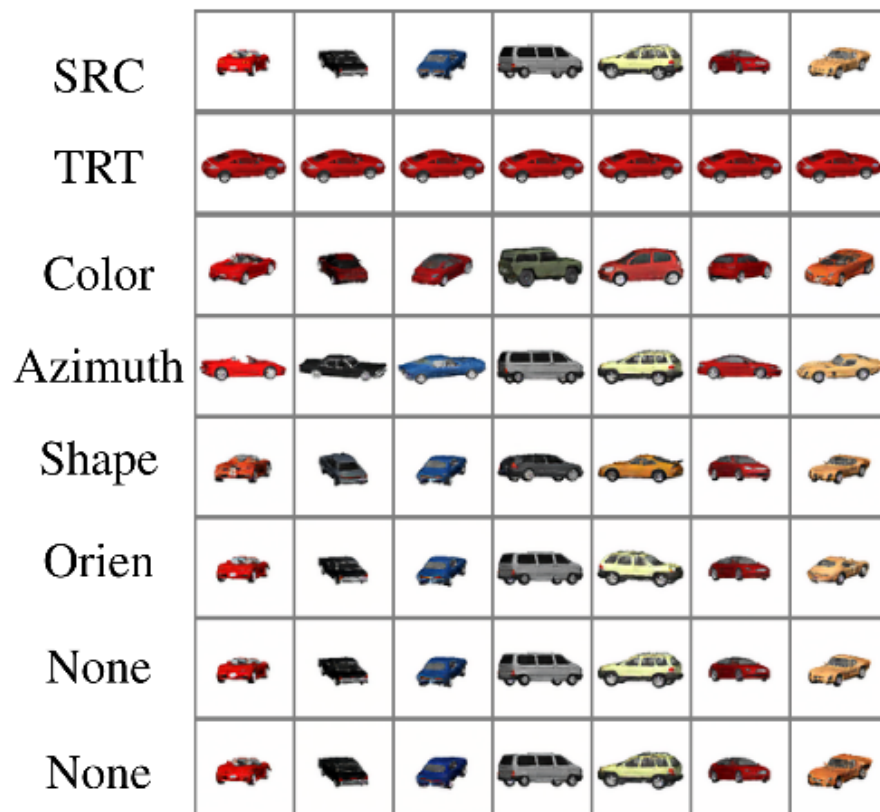


## 2. Visual Concept Learning by Disentanglement

### □ Qualitative Results



DisDiff



EncDiff (Ours)

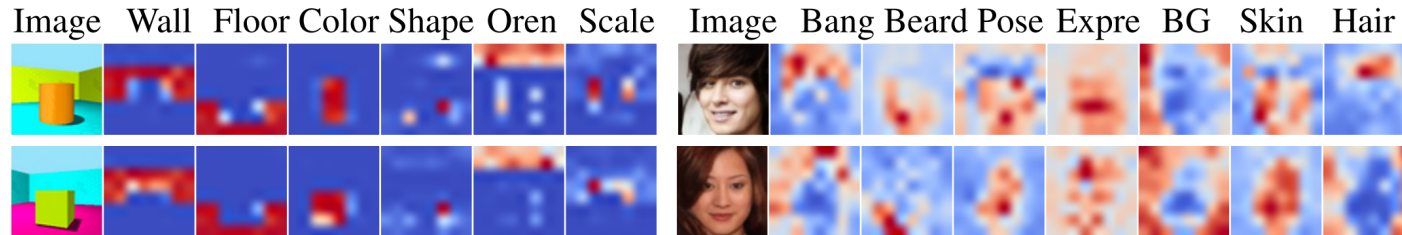
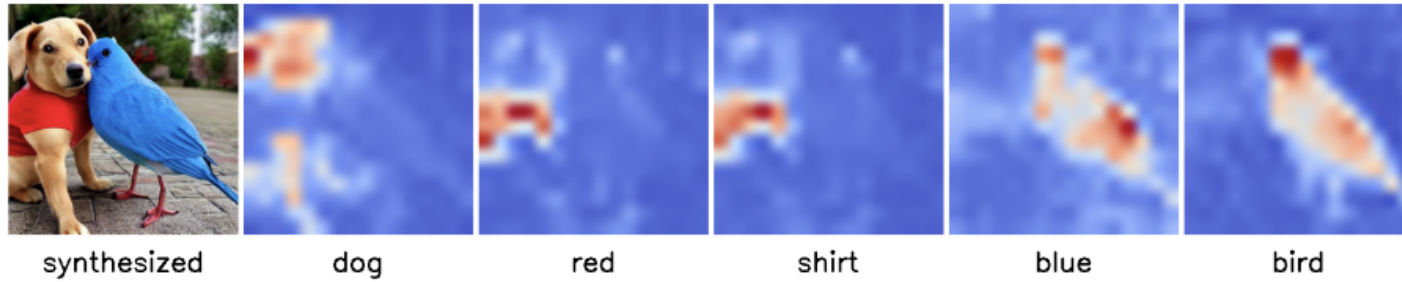
## 2. Visual Concept Learning by Disentanglement



西安交通大学  
XI'AN JIAOTONG UNIVERSITY



### □ Visualization of Attention Map





## 2. Visual Concept Learning by Disentanglement



### □ Ablation Study

Method	FactorVAE score $\uparrow$	DCI $\uparrow$
EncDiff w/sqrt	$0.997 \pm 0.011$	$0.950 \pm 0.041$
EncDiff w/sqrt linear	$0.988 \pm 0.026$	$0.924 \pm 0.050$
EncDiff w/linear	$0.999 \pm 0.002$	$0.930 \pm 0.045$
EncDiff w/cosine	<b><math>0.999 \pm 0.001</math></b>	<b><math>0.969 \pm 0.030</math></b>

Method	Params. $\downarrow$ (M)	FLOPs $\downarrow$ (M)	Time $\downarrow$ (s)
Diff-AE [24]	67.8	3955.1	31.0
DisDiff [37]	57.1	5815.8	35.3
EncDiff	42.3	2898.5	11.8

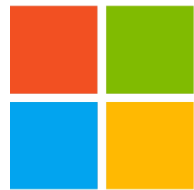
Method	FactorVAE score $\uparrow$	DCI $\uparrow$
EncDec w/o Diff	$0.537 \pm 0.074$	$0.178 \pm 0.050$
EncDiff w/ AdaGN	$0.911 \pm 0.101$	$0.637 \pm 0.068$
EncDiff	<b><math>0.999 \pm 0.000</math></b>	<b><math>0.969 \pm 0.030</math></b>

Method	FactorVAE score $\uparrow$	DCI $\uparrow$
EncDiff-V	$0.999 \pm 0.000$	$0.900 \pm 0.045$
EncDiff	<b><math>0.999 \pm 0.001</math></b>	<b><math>0.969 \pm 0.030</math></b>



西安交通大学

XI'AN JIAOTONG UNIVERSITY



Microsoft

Thank you!