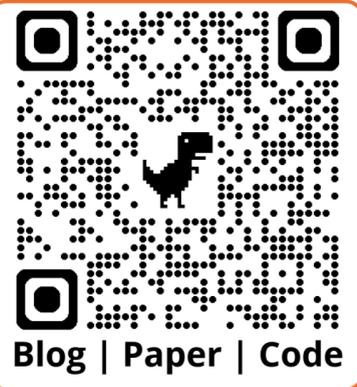


Improving Compositional Generation with Diffusion Models

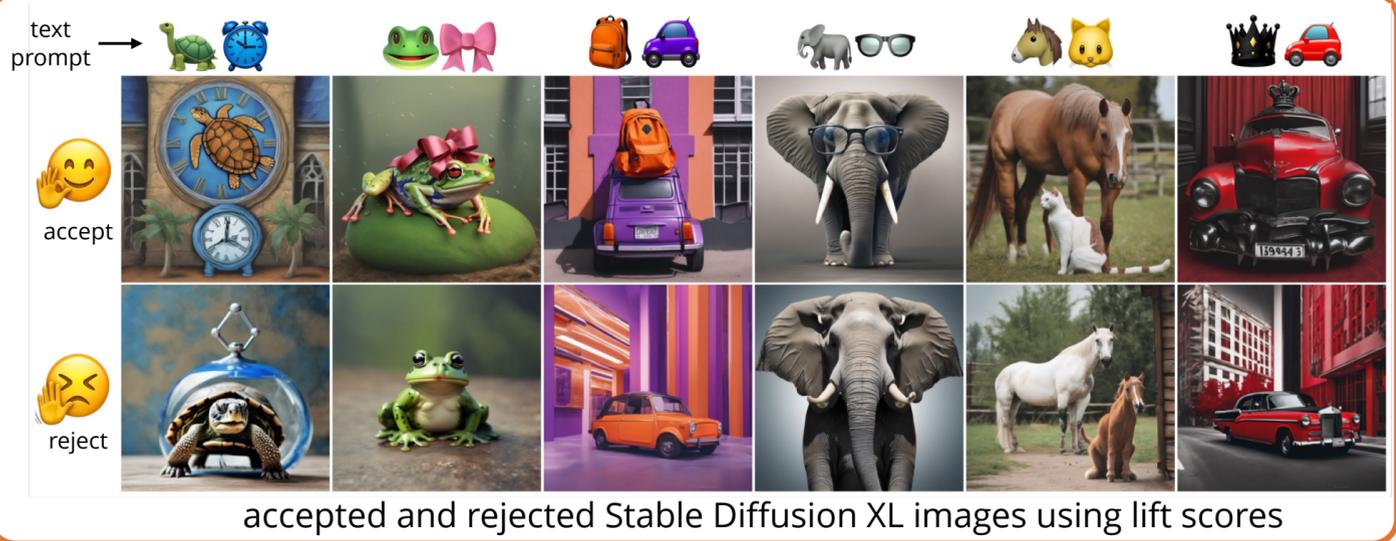
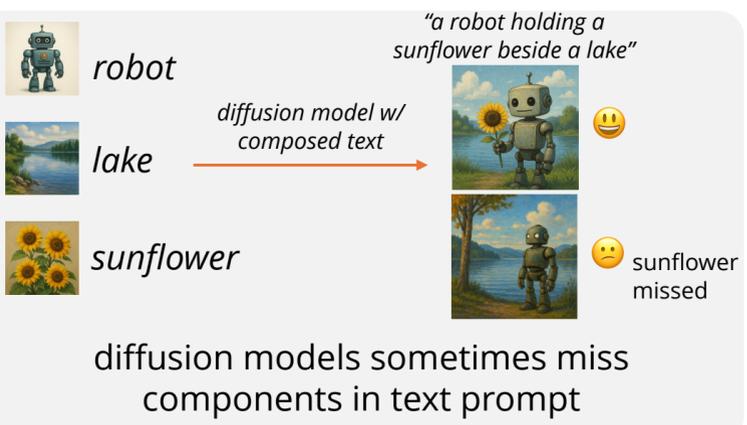
Using Lift Scores



Chenning Yu and Sicun Gao, UC San Diego

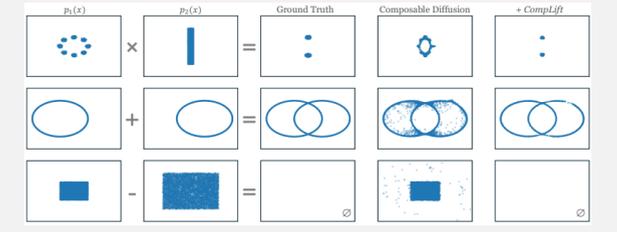
We use diffusion model itself to reject bad images, in compositional way

Background



Key Highlights

- Training-free, use diffusion model itself
- Simple concept - compare with 0 to accept/reject
- Boosts base models by 106% (2D generation) and 16% (Image generation with SD XL)



Lift Score

An existing concept in data mining (Brin et al., 1997)

$$lift(x|c) := \log \frac{p(x|c)}{p(x)}$$

sample constraint (e.g., a text component) (e.g., an image)

← conditional distribution
← unconditional distribution

Diffusion Models can Easily Estimate Lift Score

$$\log \frac{p(x|c)}{p(x)} \approx \exp(-\mathbb{E}_{t,\epsilon} \|\epsilon - \epsilon_\theta(x_t, c)\|^2 + C)$$

$$\log \frac{p(x)}{p(x)} \approx \exp(-\mathbb{E}_{t,\epsilon} \|\epsilon - \epsilon_\theta(x_t, \emptyset)\|^2 + C)$$

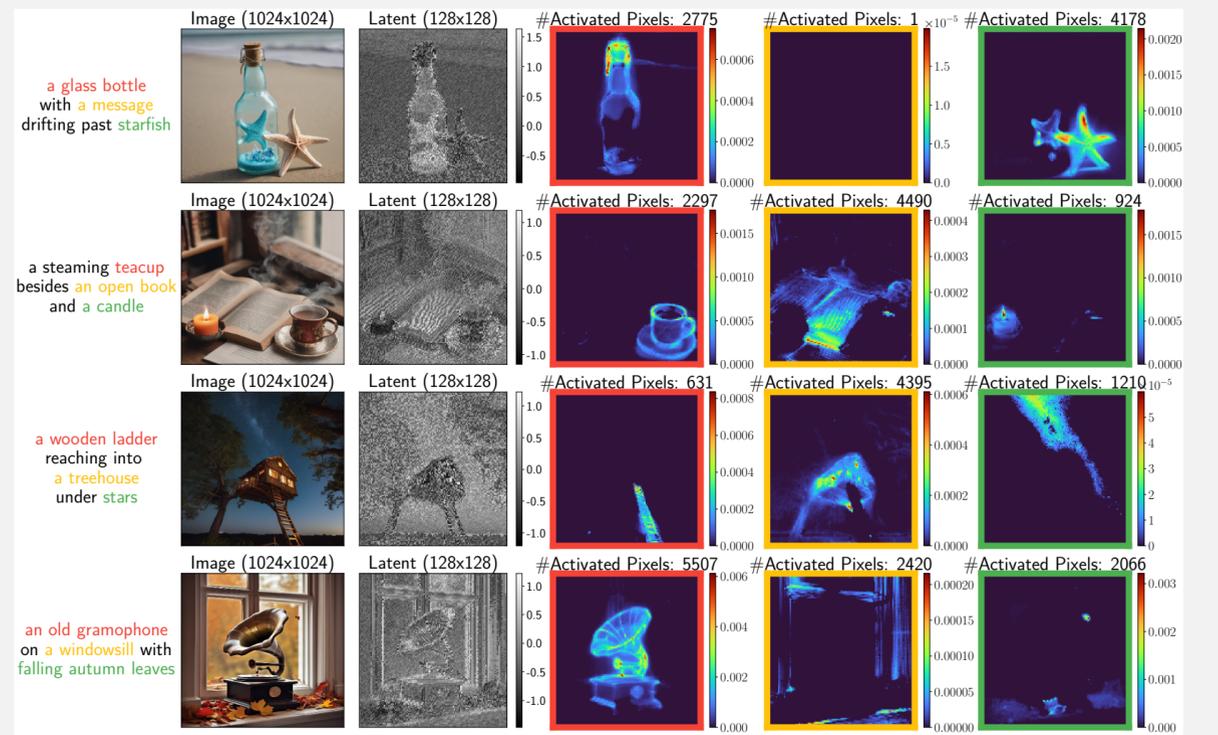
random sample timestep/noise diffusion model

We use lift score as rejection criterion to check alignment of each individual component; then compose the criterion to accept/reject

Type	Algebra	Acceptance Criterion
Product	$c_1 \wedge c_2$	$\min_{i \in [1,2]} lift(x c_i) > 0$
Mixture	$c_1 \vee c_2$	$\max_{i \in [1,2]} lift(x c_i) > 0$
Negation	$\neg c_1$	$lift(x c_1) \leq 0$

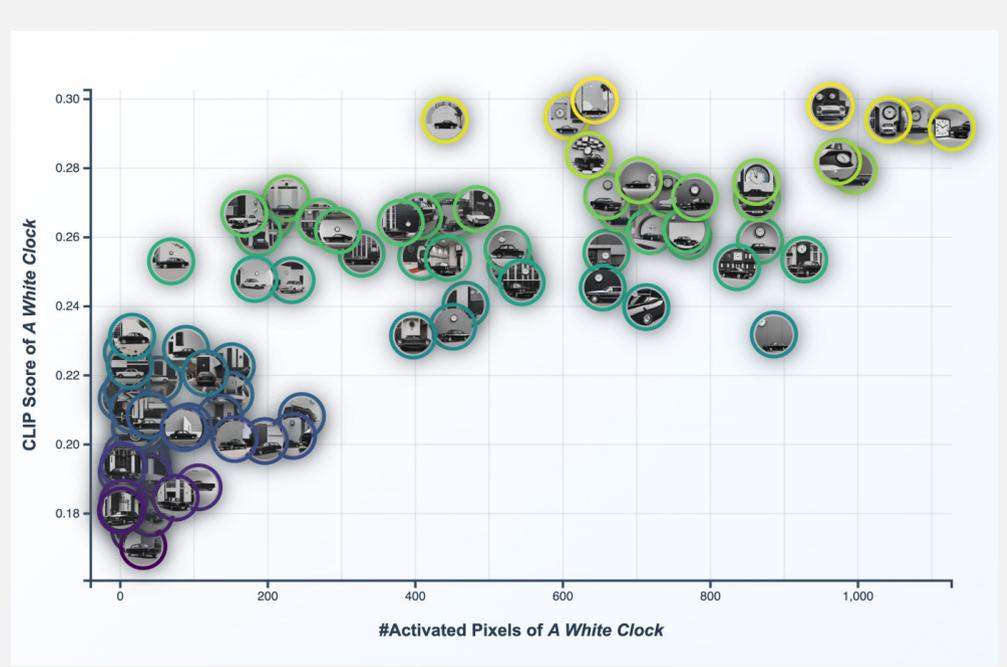
Table 1. Examples of Compose for multiple conditions.

Visualization in Latent Space



lift score shows strong correspondence to each text component; activated pixels are pixels with positive lift scores

Correlation with CLIP Score



images with missing components tend to have (1) lower CLIP score and (2) fewer activated pixels with positive lift scores
prompt: a black car and a white clock