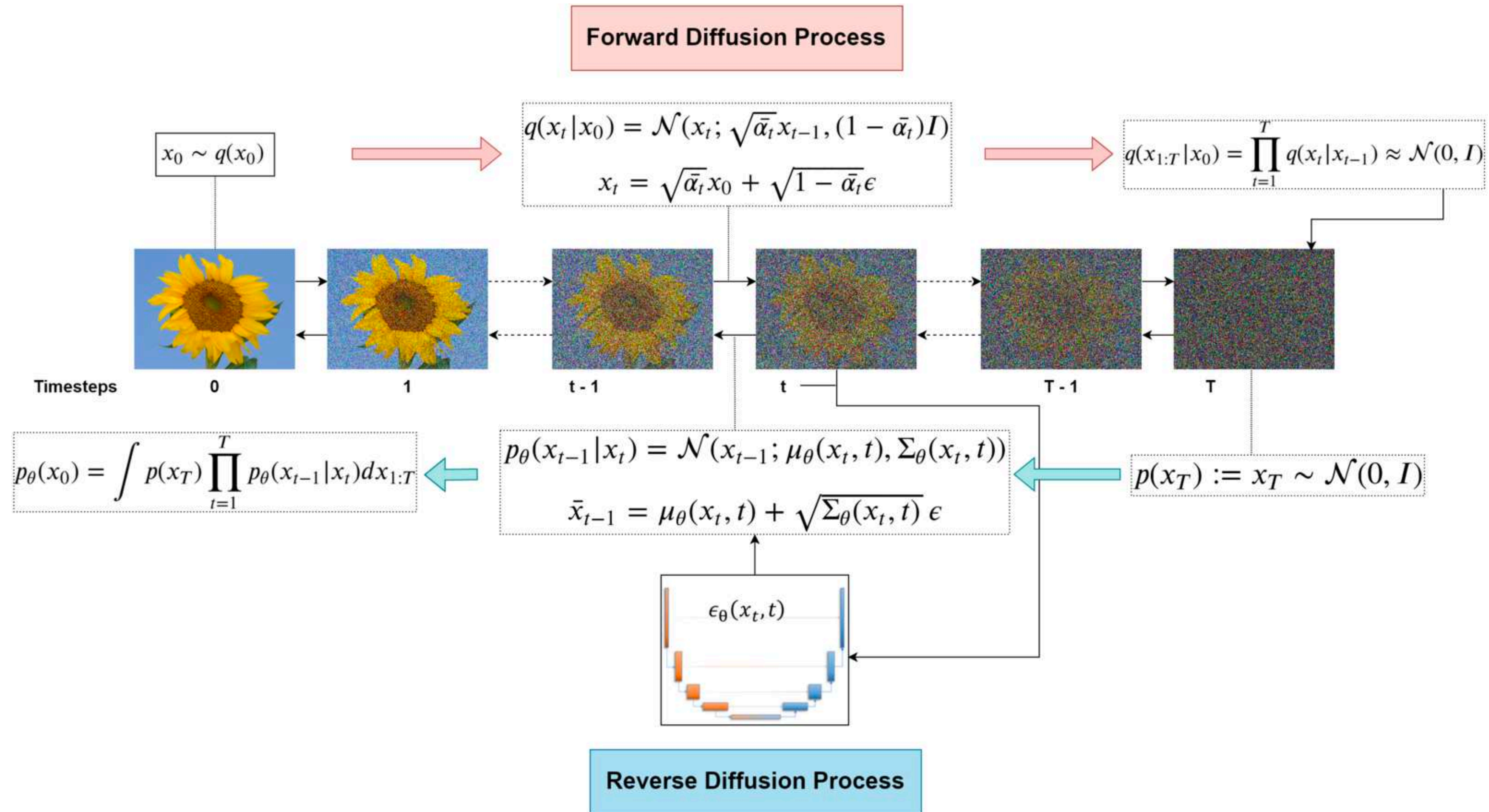


# **Diffusion Self-Guidance for Controllable Image Generation**

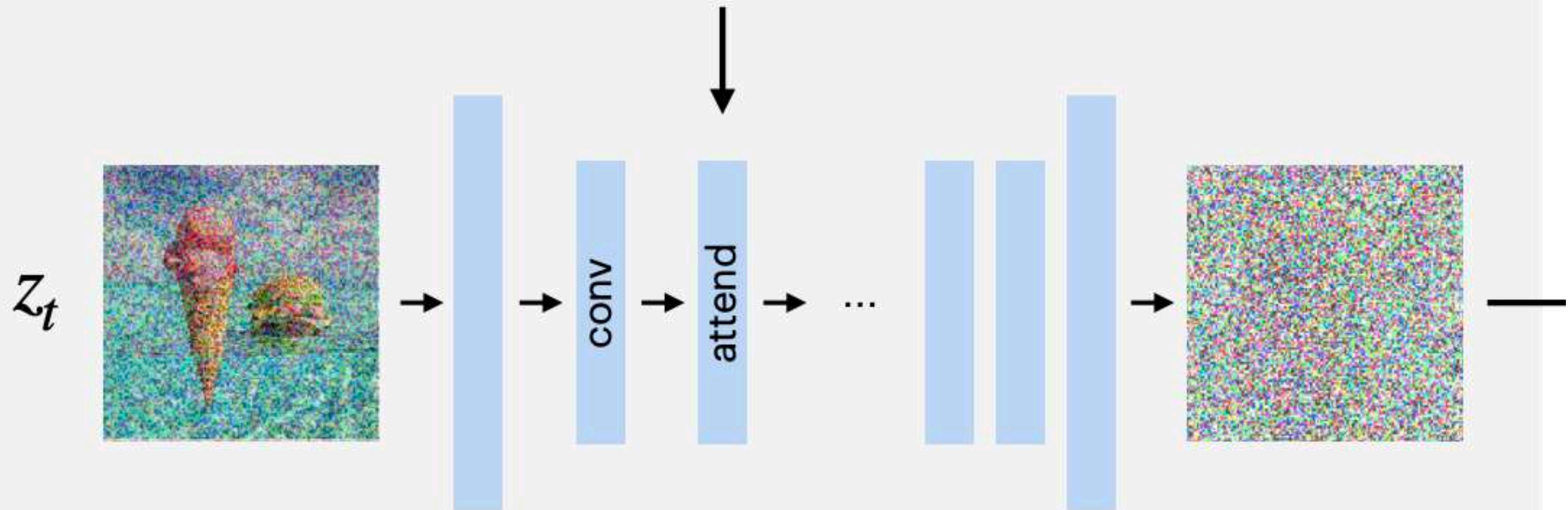
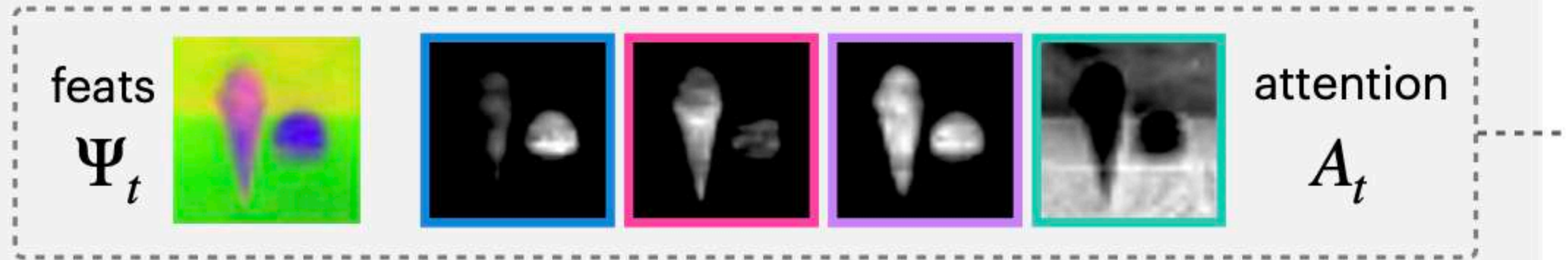
**Кокоева Мария**

# Диффузионные модели





“a photo of a burger and an ice cream cone floating in the ocean”



# Guidance

$$\hat{\epsilon}_t = \epsilon_{\theta}(z_t; t, y) - s\sigma_t \nabla_{z_t} \log p(y|z_t)$$

$$\hat{\epsilon}_t = (1 + s)\epsilon_{\theta}(z_t; t, y) - s\epsilon_{\theta}(z_t; t, \emptyset) + v\sigma_t \nabla_{z_t} g(z_t; t, y)$$



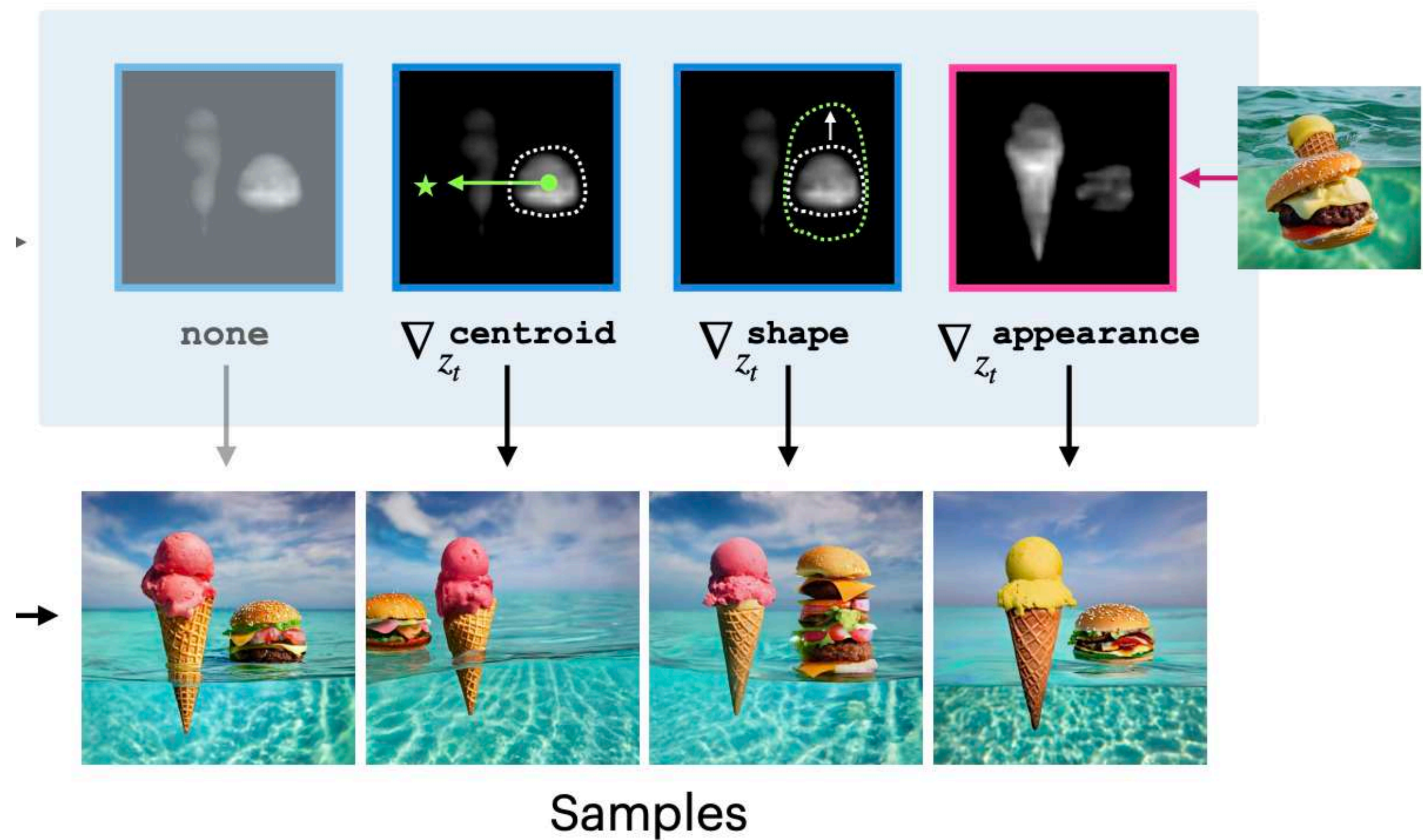
$$\text{centroid}(k) = \frac{1}{\sum_{h,w} \mathcal{A}_{h,w,k}} \begin{bmatrix} \sum_{h,w} w \cdot \mathcal{A}_{h,w,k} \\ \sum_{h,w} h \cdot \mathcal{A}_{h,w,k} \end{bmatrix}$$

$$\text{size}(k) = \frac{1}{HW} \sum_{h,w} \mathcal{A}_{h,w,k}$$

$$\text{shape}(k) = \mathcal{A}_k^{\text{thresh}}$$

$$\text{appearance}(k) = \frac{\sum_{h,w} \text{shape}(k) \odot \Psi}{\sum_{h,w} \text{shape}(k)}$$

## Self-Guidance





$$\begin{aligned}
& \overbrace{\left[ w_0 \frac{1}{|O|-1} \sum_{o \neq o_k \in O} \frac{1}{|\mathcal{A}|} \sum_{i=0}^{|\mathcal{A}|} \|\mathbf{shape}_{i,t,\text{orig}}(o) - \mathbf{shape}_{i,t}(o)\|_1 \right]}^{\text{Fix all other object shapes}} \\
& + \overbrace{w_1 \frac{1}{|O|} \sum_{o \in O} \|\mathbf{appearance}_{t,\text{orig}}(o) - \mathbf{appearance}_t(o)\|_1}^{\text{Fix all appearances}} \\
& + \overbrace{w_2 \frac{1}{|\mathcal{A}|} \sum_{i=0}^{|\mathcal{A}|} \|\mathcal{T}(\mathbf{shape}_{i,t,\text{orig}}(o_k)) - \mathbf{shape}_{i,t}(o_k)\|_1}^{\text{Guide } o_k \text{'s shape to translated original shape}}
\end{aligned}$$

“distant shot of the tokyo tower with a massive sun in the sky”



“a photo of a fluffy cat sitting on a museum bench looking at an oil painting of cheese”



“a photo of a raccoon in a barrel going down a waterfall”



(a) Original   (b) Move up   (c) Move down   (d) Move left   (e) Move right   (f) Shrink   (g) Enlarge



Fix all object shapes

$$g = w_0 \frac{1}{|O|} \sum_{o \in O} \frac{1}{|\mathcal{A}|} \sum_{i=0}^{|\mathcal{A}|} \|\text{shape}_{i,t,\text{orig}}(o) - \text{shape}_{i,t}(o)\|_1$$

“a photo of a parrot riding a horse down a city street”



“a photo of a bear wearing a suit eating his birthday cake out of the fridge in a dark kitchen”



(a) Original

(b) New appearances

(c) ControlNet [39]

(d) PtP [11]



$$\begin{aligned}
g = & \underbrace{w_0 \frac{1}{|O|} \sum_{o \in O} \frac{1}{|\mathcal{A}|} \sum_{i=0}^{|\mathcal{A}|} \|\text{shape}_{i,t,A}(o) - \text{shape}_{i,t}(o)\|_1}_{\text{Copy object shapes from A}} \\
& + \underbrace{w_1 \frac{1}{|O|} \sum_{o \in O} \|\text{appearance}_{t,B}(o) - \text{appearance}_t(o)\|_1}_{\text{Copy object appearance from B}}
\end{aligned}$$





$$g = w_0 \underbrace{\frac{1}{J} \sum_j \frac{1}{|\mathcal{A}|} \sum_{i=0}^{|\mathcal{A}|} \|\text{shape}_{i,t,j}(o_{k_j}) - \text{shape}_{i,t}(o_k)\|_1}_{\text{Copy each object's shape, position, and size}} + w_1 \underbrace{\frac{1}{J} \sum_j \|\text{appearance}_{t,j}(o_{k_j}) - \text{appearance}_t(o_k)\|_1}_{\text{Copy each object's appearance}}$$

“a photo of a picnic blanket, a fruit tree, and a car by the lake”



(a) Take blanket



(b) Take tree



(c) Take car



(d) **Result**



(e) + Target layout



(f) **Final result**

“a top-down photo of a tea kettle, a bowl of fruit, and a cup of matcha”



(a) Take matcha



(b) Take kettle



(c) Take fruit



(d) **Result**



(e) + Target layout



(f) **Final result**

“a photo of a dog wearing a knit sweater and a baseball cap drinking a cocktail”



(a) Take sweater



(b) Take cocktail



(c) Take cap



(d) **Result\***

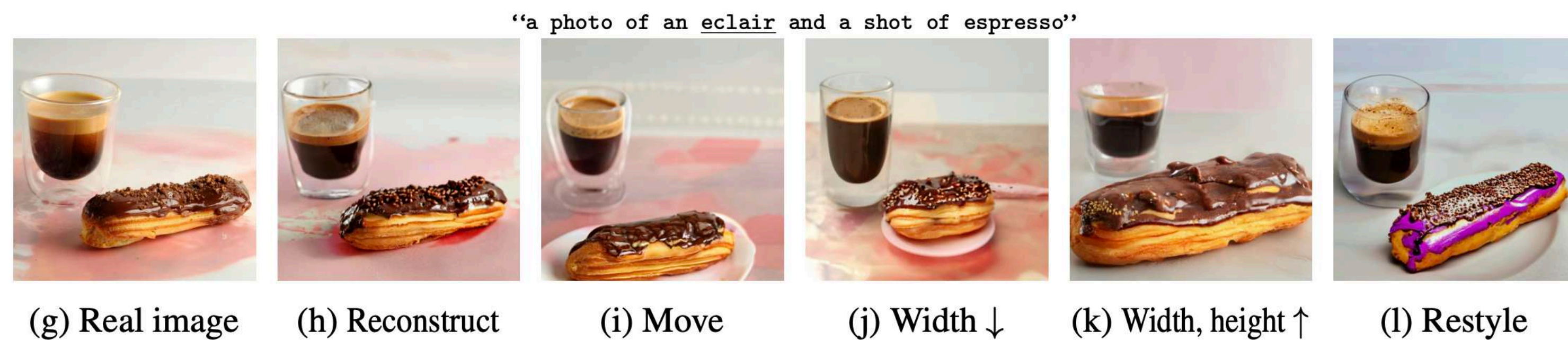
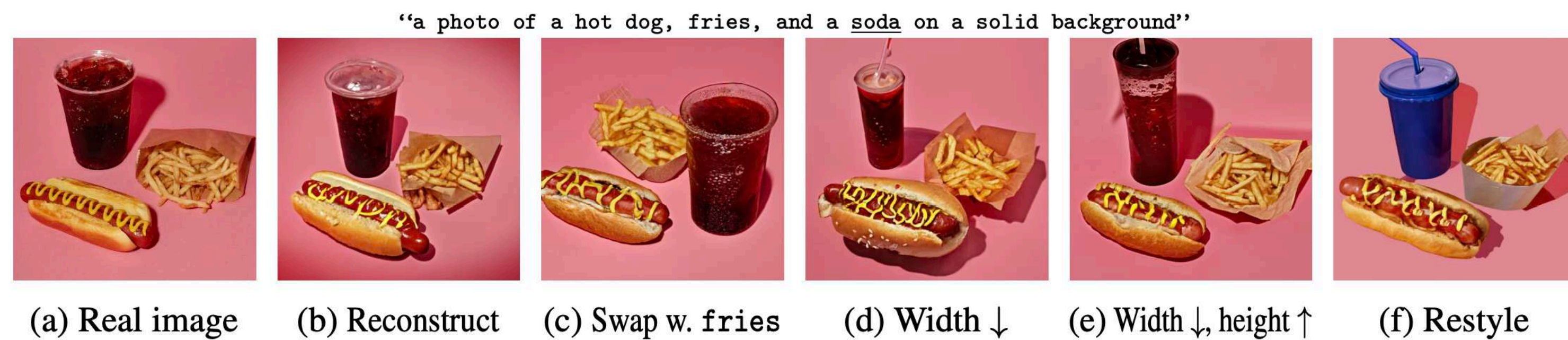


(e) + Target layout



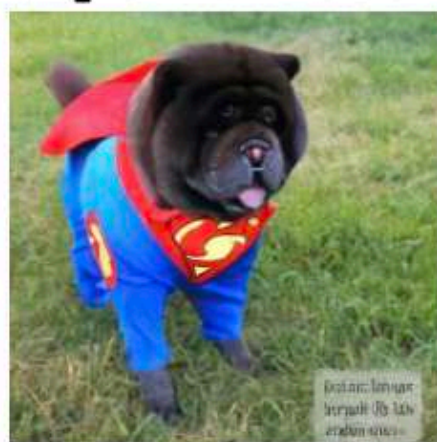
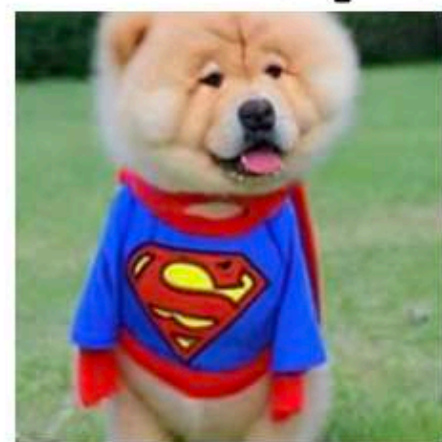
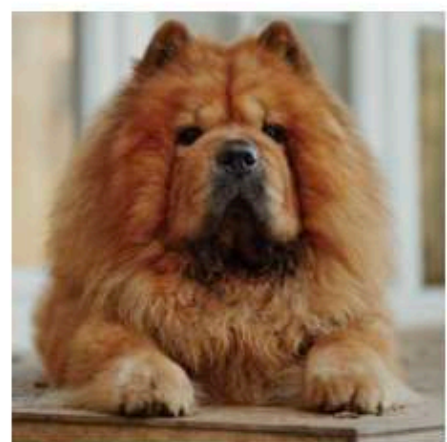
(f) **Final result**



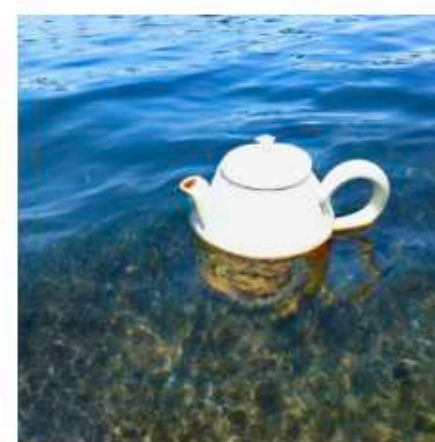
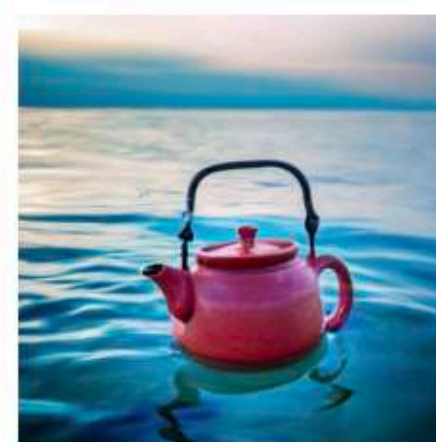
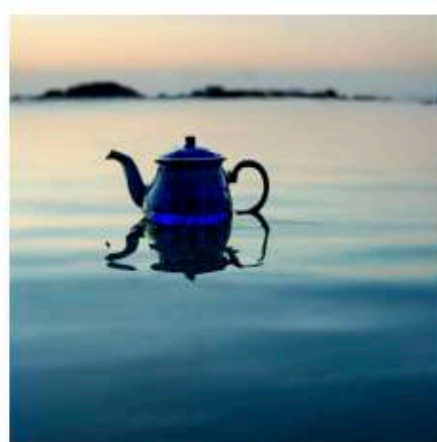
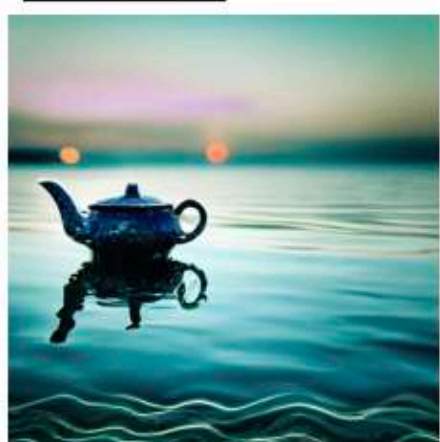
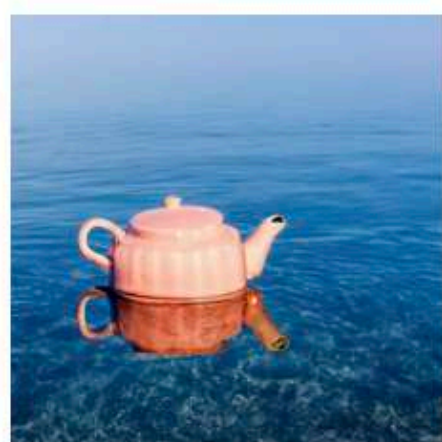
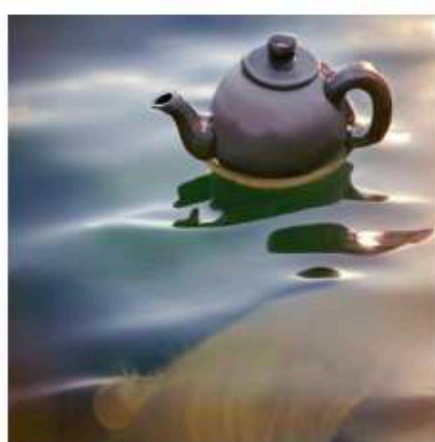




“a photo of a chow chow wearing a superman outfit”



“a dslr photo of a teapot floating in the sea”



(a) Original

(b) Ours

(b) Random samples without self-guidance