Original: "a photo of a giant macaron and a croissant in the seine with the eiffel tower visible"

Original: "a photo of a meatball and a donut falling from the clouds onto a neighborhood"



Figure 1: Hyperparameter sweeps. We show results for two edits (moving and resizing objects) using Eqn. 4 in the Supp. Mat., for different values of weights on the three edit terms, holding the other terms to the value in the middle column. Please zoom in to view in more detail. Reasonable values in the middle columns (within the expected range) lead to overall successful image manipulation. Very large hyperparameter values cause visual artifacts to appear (by moving sampling off-manifold) while still tending to perform the edit successfully, while extremely small values often fail to conduct the edit, inducing artifacts resulting from a "half-executed" manipulation.

"a watermelon and a pitcher of beer on a picnic table"

"a sea otter playing volleyball at the beach"





frog centroid (0.3, 0.5)



frog centroid (0.65, 0.5)

"a hot air balloon and a flock of geese flying through the sky on top of new york"



balloon size 0.05



balloon size 0.2

globe centroid (0.35, 0.35)



cat size 0.03



globe centroid (0.9, 0.9)

cat size 0.25

Figure 2: Self-guidance on Stable Diffusion XL. To highlight the generality of our approach, we demonstrate preliminary results for controllable generation on a popular latent-space text-to-image diffusion model, using 100 DDPM steps (applying self-guidance from step 10 to step 90). We get best results only guiding attention in the second decoder block of the denoiser model.