# REMinD: Balancing Robust Concept Unlearning and Image Quality in Diffusion Models

## Supplementary Material

## A. ADDITIONAL EXPERIMENTS

### A.1. Additional Visualizations Results for Toxic Unlearning

In Section **Toxic Unlearning Performance**, we compare the performance of unlearning the concept 'nudity'. However, inappropriate prompts also cover other categories. Fig. A1 shows the performance of three generative models (SD-v1.4, ESD, and REMinD) in handling prompts related to different toxic categories, such as shocking, self-harm, violence, and harassment. SD-v1.4 serves as the baseline model, while ESD and REMinD have unlearned inappropriate concepts, avoiding generating related concepts when handling sensitive content. The results of the REMinD model are of significant importance for the societal application of generative models, especially in scenarios that require strict content control.

### A.2. Additional Visualizations Results for Object Unlearning

In Section **Robustness against Adversarial Prompt Attacks**, we demonstrate the effectiveness of our defense against adversarial attacks. Here, we also present additional generated results from the models. Fig. A2 highlights the performance of REMinD (configured with various unlearning targets such as Church, Garbage Truck, Parachute, and Tench) in generating images across four global categories: kitchen, skiing, cat, and toilet. The prompts used include a variety of descriptive scenarios, such as a gas stove in a kitchen, a black-and-white image of a man skiing, a kitten sitting in a dish, and a modern toilet in a tiled bathroom. REMinD excels in object unlearning, showing resilience to adversarial attacks while keeping generative capabilities similar to the original model. This balance between mitigating undesired concepts and preserving the quality of the output highlights the potential of REMinD for practical applications requiring both safety and high performance in image generation.

### A.3. Extra Visualizations Results on UnlearnCanvas for Unlearning

We provide additional visualization results on the UnlearnCanvas dataset to further demonstrate the effectiveness of our method in concept unlearning. Figures A3 and A4 illustrate the results of **object unlearning** under fixed *style* conditions. Specifically, Figure A3 shows results in the *Photo* style, while Figure A4 presents the corresponding results in



Figure A1. More visualization results of toxic unlearning. Each column is a different prompt in a different category. Similar to 'nudity', these categories are another toxic type of prompts.



Figure A2. More visualization results of object unlearning. We use our REMinD to unlearn different hypernym concepts: 'church', 'garbage truck', 'parachute', and 'tench', Then, we compare the performance of generation ability on global concepts, such as 'kitchen', 'skiing', 'cat', and 'toilet'.

the *Van Gogh* style. Each column corresponds to a different object concept, and the diagonal images represent the targets that have been unlearned. The reduction of visual

fidelity for the erased concepts, alongside the preservation of unrelated concepts, indicates successful object-specific unlearning.

In contrast, Figures A5 and A6 focus on **style unlearning** under fixed *object* conditions. Figure A5 demonstrates style unlearning for the *Dogs* object, and Figure A6 shows the results for the *Waterfalls* object. Each column represents a different artistic style, with diagonal samples indicating the target styles to be unlearned. These results show that REMinD effectively removes undesired style concepts while maintaining coherent generation quality across other styles, confirming the flexibility and robustness of our approach to both object and style unlearning scenarios.
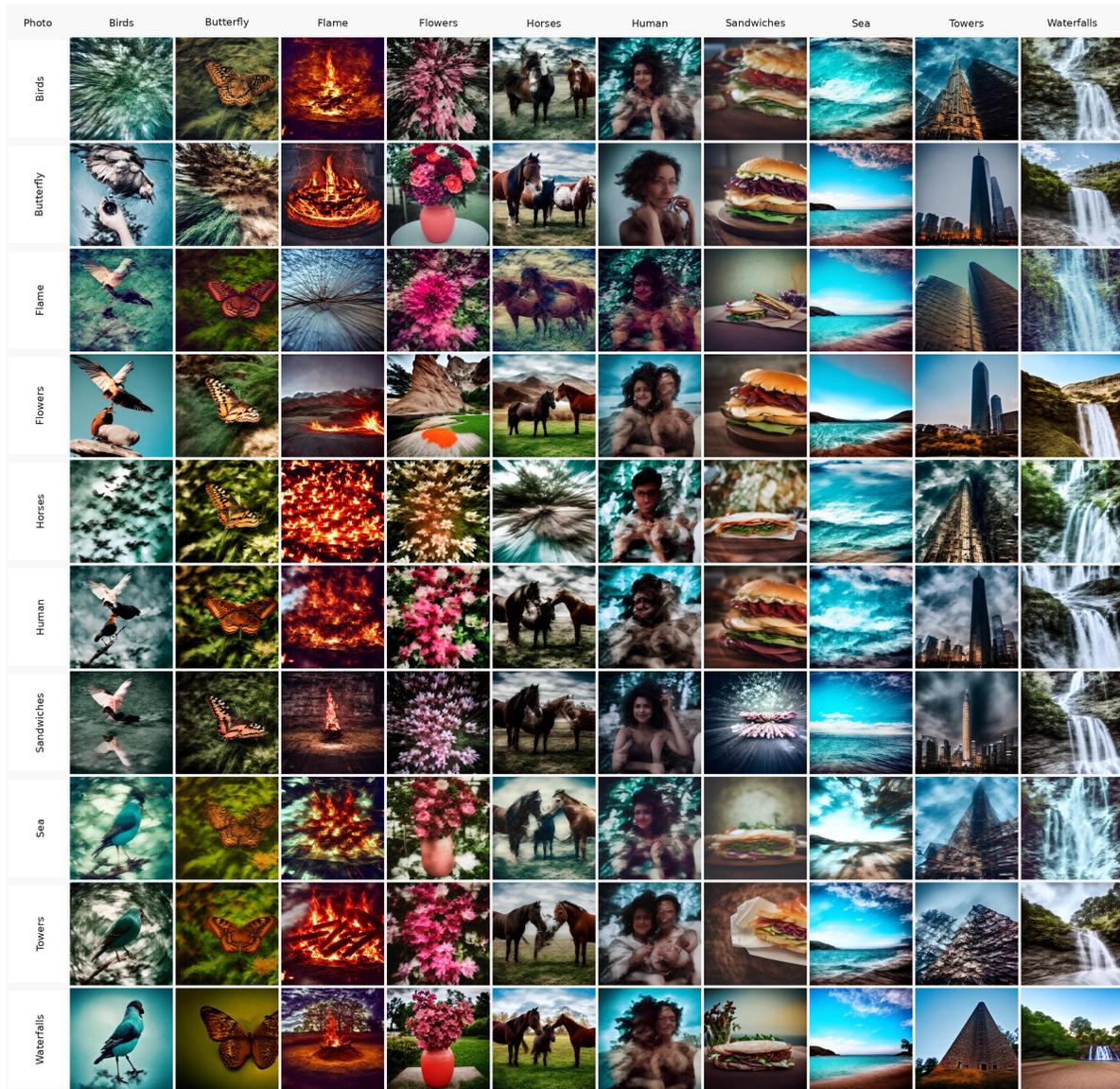
Figure A3. Object unlearning on UnlearnCanvas Dataset. All images are set in Photo style. Each column is a different concept. The images on the diagonal are the targets of unlearning.
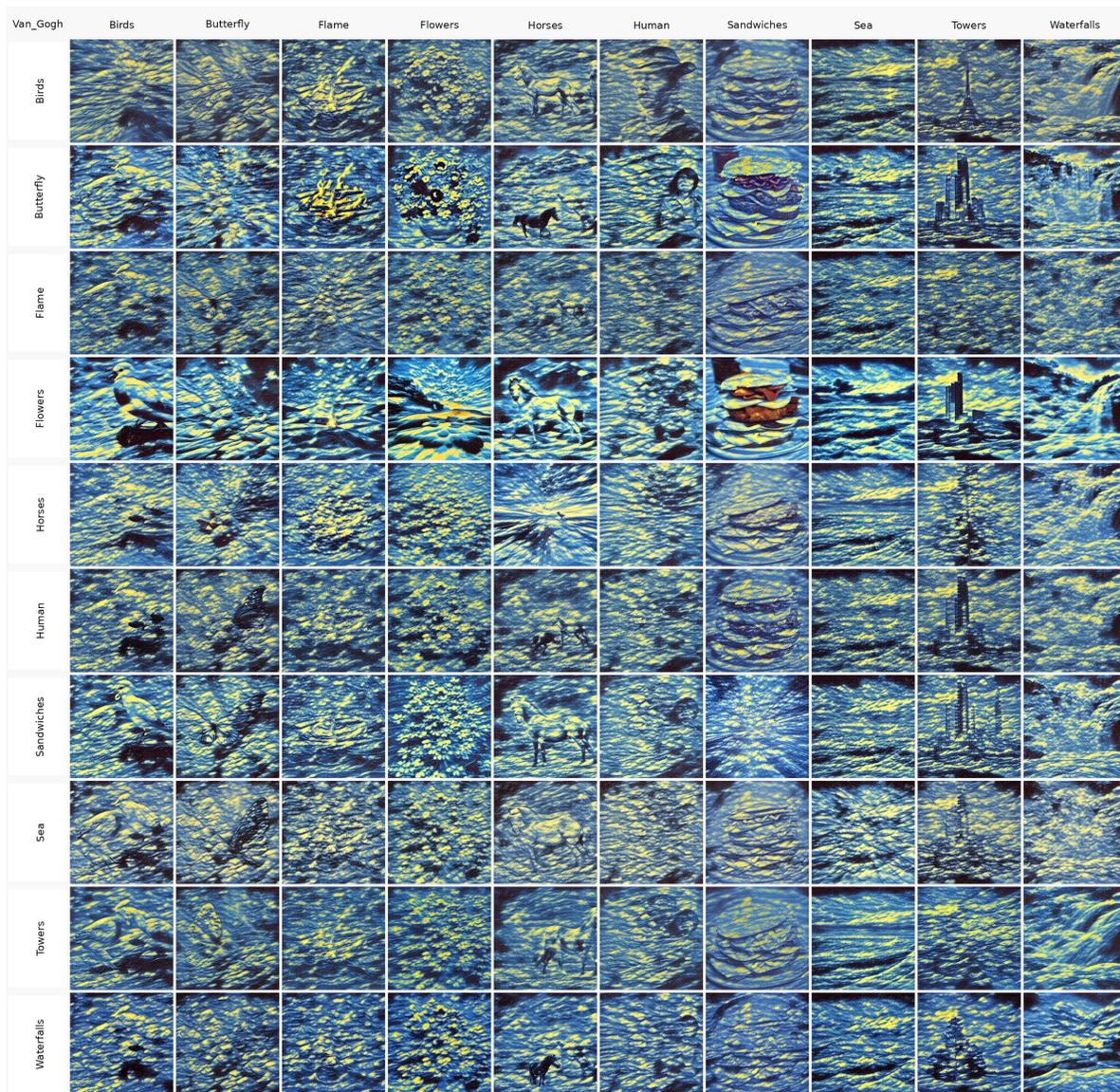
Figure A4. Object unlearning on UnlearnCanvas Dataset. All images are set in Van Gogh style. Each column is a different concept. The images on the diagonal are the targets of unlearning.
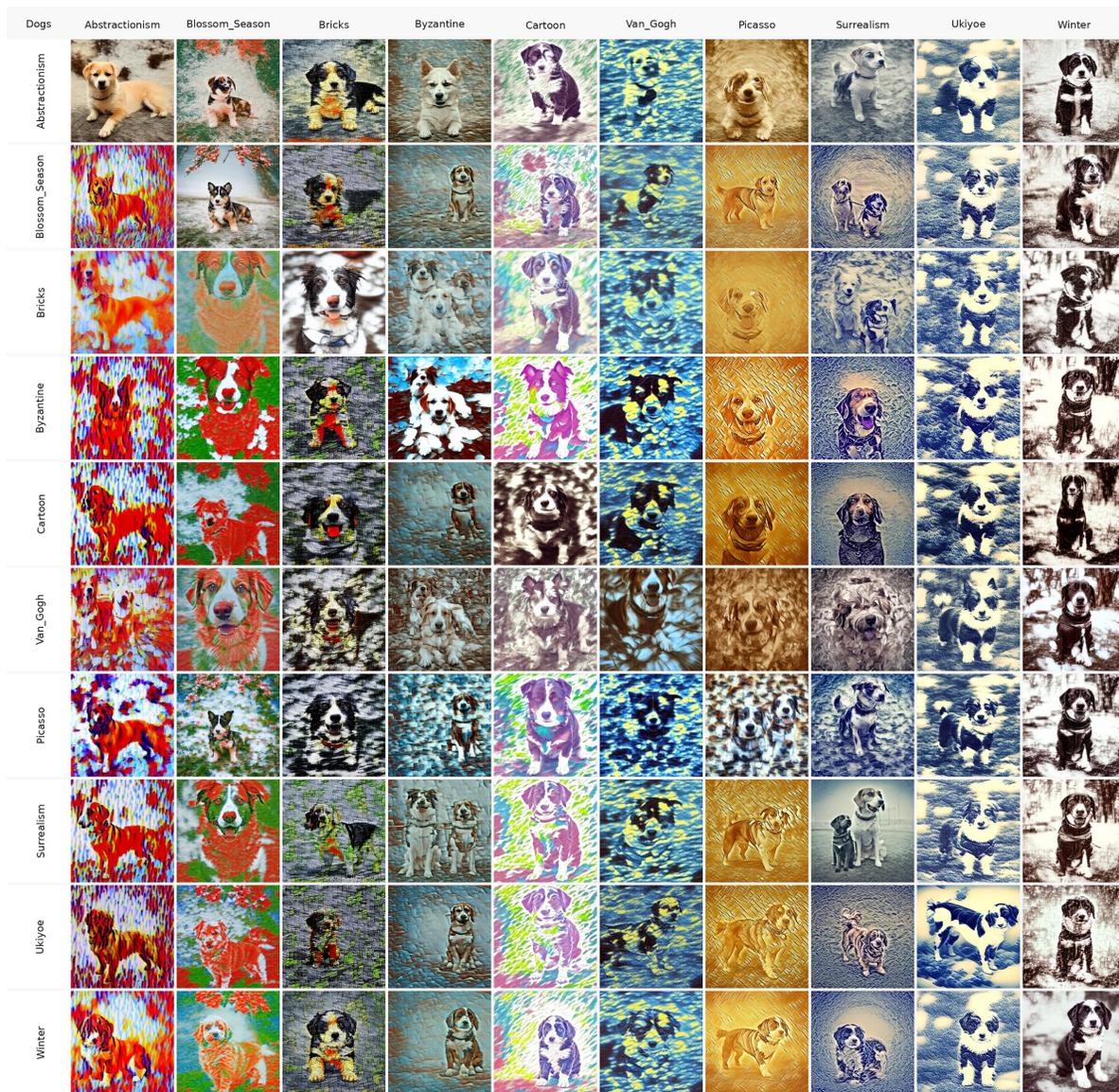
Figure A5. Style unlearning on UnlearnCanvas Dataset. All images are set in Dogs object. Each column is a different style. The images on the diagonal are the targets of unlearning.
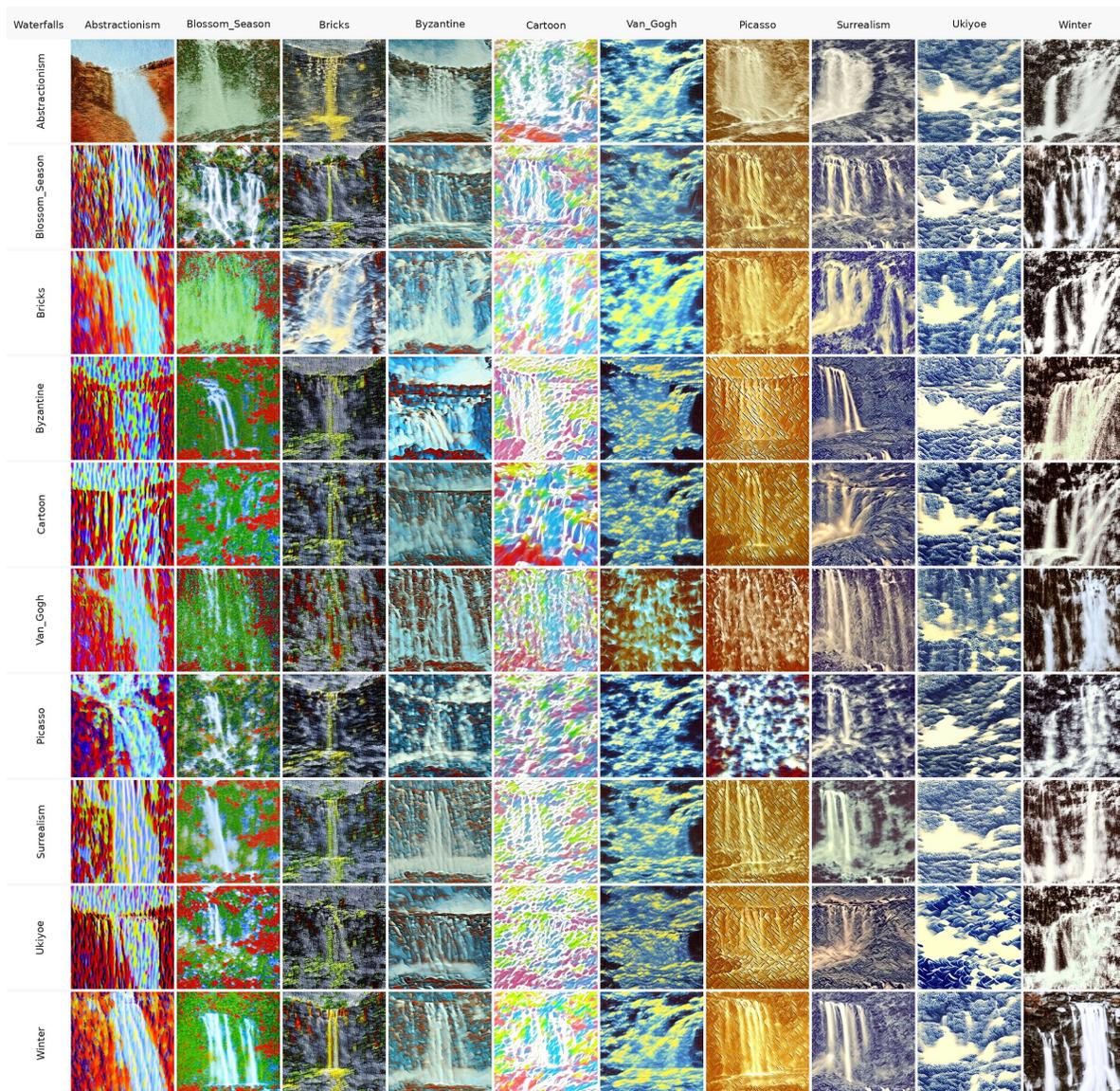
Figure A6. Style unlearning on UnlearnCanvas Dataset. All images are set in Waterfalls object. Each column is a different style. The images on the diagonal are the targets of unlearning.