

Memorization In Stable Diffusion Is Unexpectedly Driven by CLIP Embeddings

Supplementary Material

A. Additional Experimental Results

A.1. Applying Our Method To Stable Diffusion v1.4

We apply our mitigation method, consisting of $\langle \text{pad} \rangle$ replacement and v^{eot} masking. As shown in Figure 9, our method significantly reduces memorization without degrading image quality.

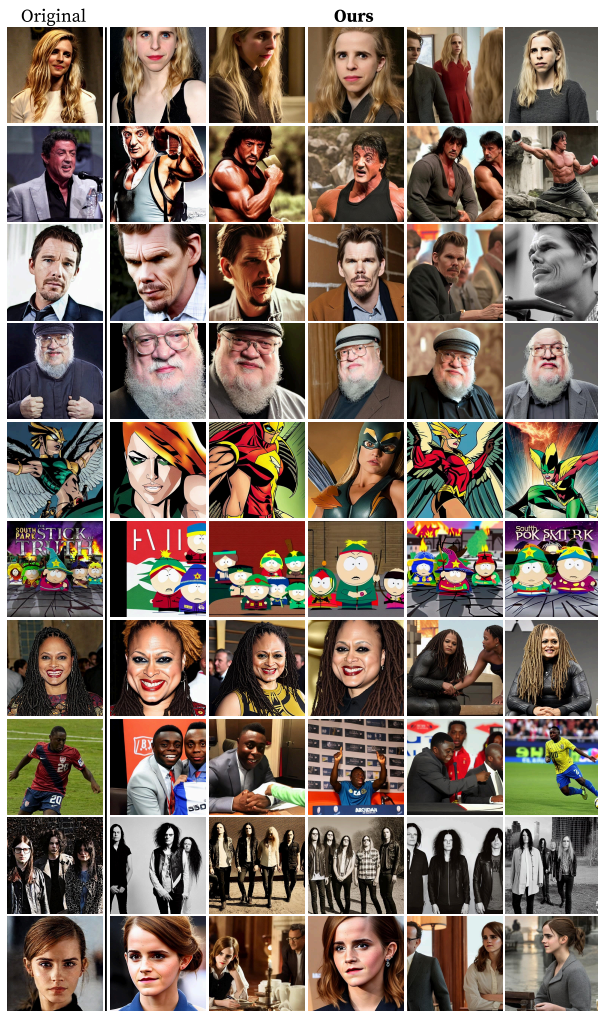


Figure 9. $\langle \text{pad} \rangle$ replacement and v^{eot} masking. The first column (Original) shows the memorized image consistently reproduced from the original embedding regardless of seeds. The remaining five columns (Ours) are generated using our mitigation method with five different random seeds.

A.2. Comparison with Prior Mitigation Methods

Comparison with Baselines (1) Ren et al. [24], which rescales cross-attention. (2) Wen et al. [36], which minimize the magnitude of text-conditional noise prediction. (3) Random Token Addition (RTA) [31], which inserts randomly sampled tokens. (4) Random Number Addition (RNA) [31], which inserts a random number between $[0, 10^6]$. We use their optimal configurations $\text{optim_target_loss} = 3$ for Wen et al. [36] and $\text{rescale_attention} = 1.25$ for Ren et al. [24]. Results are shown in Figure 10.

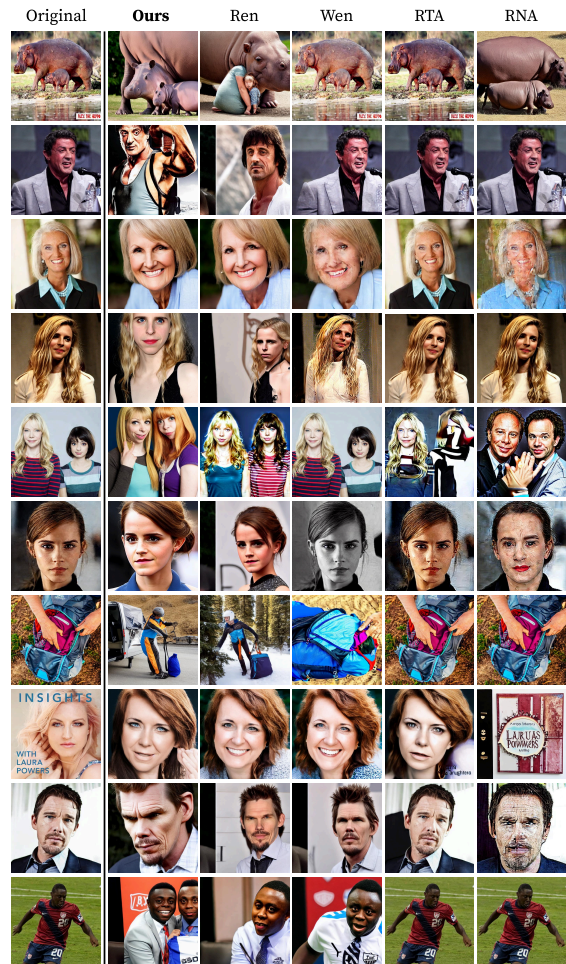


Figure 10. Each row shows images generated from the same prompt using different mitigation methods. Visually and quantitatively, our method produces a more favorable quality-preserving trade-off.

Inference Time Comparison. As shown in Table 6, our mitigation method achieves the fastest inference among all evaluated approaches. Since it does not require any additional optimization or computation during inference, its runtime closely matches that of standard text-to-image generation without mitigation. In contrast, other methods incur additional computational overhead, resulting in slower inference.

Method	Inference Time (s)
Ours	3.17 ± 0.01
Ren et al. [24]	3.37 ± 0.02
Wen et al. [36]	3.55 ± 0.16
RNA [31]	3.26 ± 0.02
RTA [31]	3.31 ± 0.02

Table 6. Average inference time per image for each method.

Evaluation on Webster 500 Prompts. Many prior memorization studies [4, 18, 24, 36] report results directly on all 500 candidate prompts released by Webster [35]. In the main paper, we use a stricter curated benchmark of 458 MV prompts, constructed from Webster [35] and Membench [13], for primary analysis, as the full 500-prompt set includes unstable and degraded cases that can confound mechanistic interpretation. At the same time, to facilitate direct comparison with this widely used evaluation setting, we additionally evaluate all mitigation methods on the full Webster 500 set.

Method	SSCD	CLIPScore	Aesthetic
Ours	0.39 ± 0.20	0.29 ± 0.03	5.33 ± 0.26
Wen	0.65 ± 0.22	0.31 ± 0.03	5.35 ± 0.23
Ren	0.20 ± 0.15	0.26 ± 0.03	4.90 ± 0.31
RTA	0.70 ± 0.23	0.31 ± 0.03	5.36 ± 0.23
RNA	0.70 ± 0.23	0.32 ± 0.03	5.33 ± 0.25

Table 7. Comparison of mitigation methods on all 500 candidate prompts from Webster [35].

As shown in Table 7, our method ($\langle \text{pad} \rangle \rightarrow !$ & $\mathbf{v}^{\text{eot}} \rightarrow \mathbf{0}$) continues to provide a favorable trade-off between memorization reduction and generation quality on this broader benchmark. In particular, compared to Wen et al. [36], RNA [31], and RTA [31], our method achieves substantially lower SSCD while maintaining comparable CLIPScore and Aesthetic Score. Although Ren et al. [24] attains lower SSCD, it does so at a noticeably larger cost in prompt alignment and visual quality. These results are consistent with our main-paper findings: even on the noisier Webster 500 set, our method remains a competitive mitigation method with a strong quality-preserving trade-off.

Evaluation on Membench. To ensure the robustness of our findings, we conduct extensive evaluations across multiple datasets. As shown in Table 8, on the Membench [13], which consists of 3,000 prompts, our method shows consistently comparable or better results against existing baselines. We do not report RNA or RTA results, as their performance is not comparable to these three representative baselines.

Method	SSCD	CLIPScore	Aesthetic
Ours	0.16 ± 0.23	0.30 ± 0.07	5.14 ± 0.53
Ren	0.15 ± 0.19	0.29 ± 0.06	5.17 ± 0.51
Wen	0.79 ± 0.37	0.31 ± 0.06	5.22 ± 0.48

Table 8. Comparison of mitigation methods on Membench [13].

A.3. Cross-Prompt Swapping Experiment

To further investigate the role of \mathbf{v}^{eot} , we conduct an experiment swapping \mathbf{v}^{eot} and $\mathbf{v}_i^{\text{pad}}$ between different prompts. Our experiment reveals that swapping only \mathbf{v}^{eot} has no discernible effect on the output, whereas swapping both \mathbf{v}^{eot} and $\mathbf{v}_i^{\text{pad}}$ alters the generated image to match the swapped prompt (Figure 11). This suggests that $\mathbf{v}_i^{\text{pad}}$, much like \mathbf{v}^{eot} , encode critical semantic information rather than serving as mere placeholders.

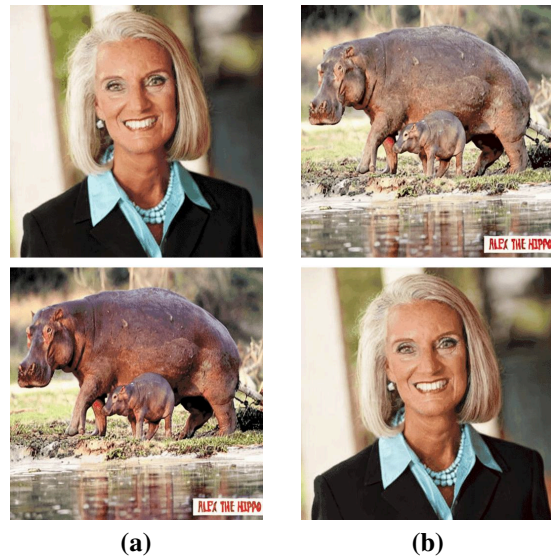
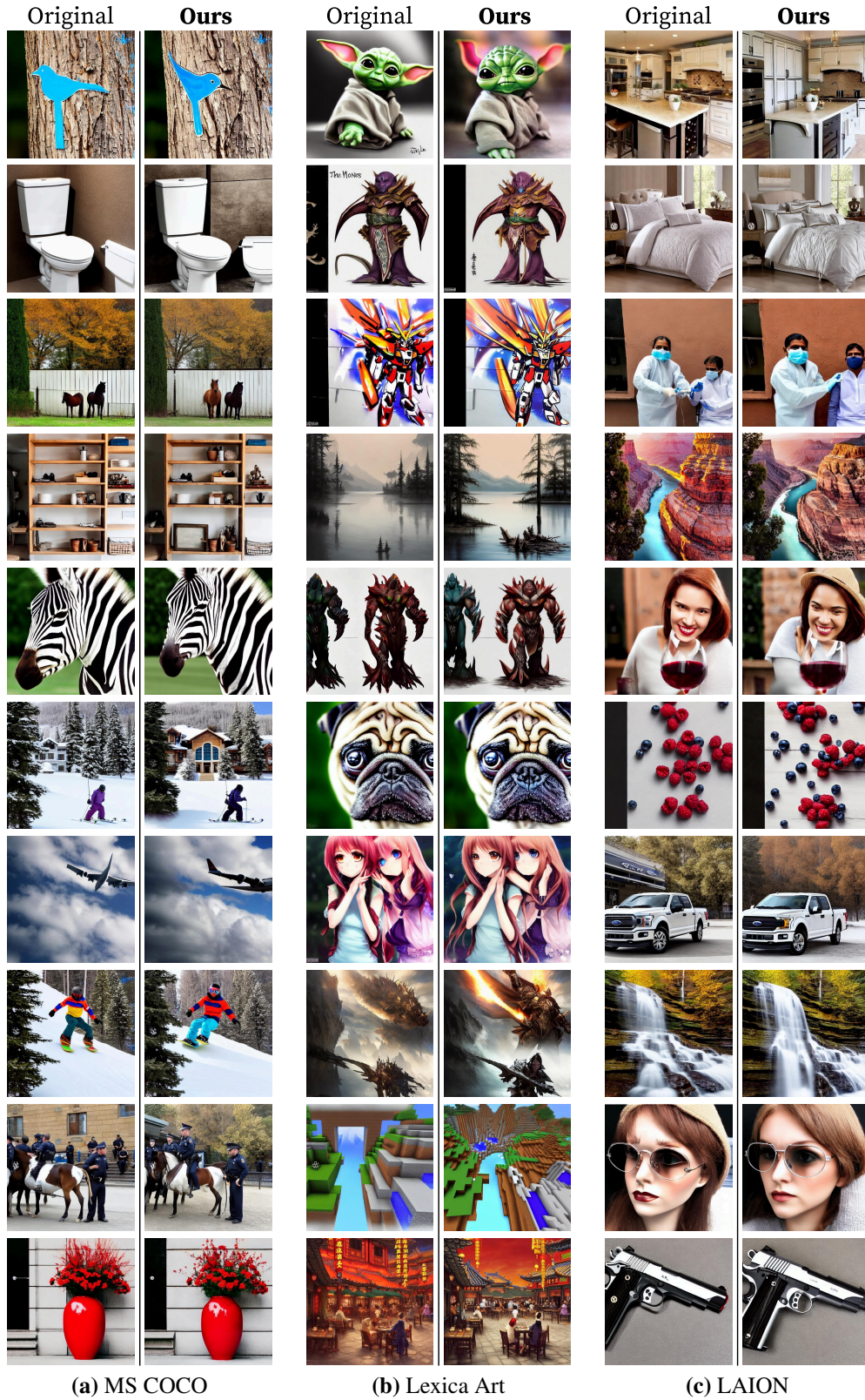


Figure 11. (a) Original generation without swapping. (b) Generation after swapping both \mathbf{v}^{eot} and $\mathbf{v}_i^{\text{pad}}$, showing that the output changes to reflect the swapped prompt.

A.4. Robustness of our mitigation method on non-memorized prompts

We utilize non-memorized prompts to evaluate whether our mitigation method impacts general generation quality. Fol-



(a) MS COCO

(b) Lexica Art

(c) LAION

Figure 12. “Original” is generated using the default embeddings, while “Ours” applies our mitigation strategy ($\langle \text{pad} \rangle \rightarrow !$ & $\mathbf{v}^{\text{eot}} \rightarrow \mathbf{0}$).

Metric	SD1		SD2
	Webster [35]	Hong et al. [13]	Webster [35]
Number of Prompts	345	3000	219
Mean Length	19.74	15.31	21.51
Max / Min Length	38 / 4	38 / 5	35 / 6

Table 9. Prompt Token Length Statistics for Stable Diffusion v1.4 (SD1) and Stable Diffusion v2.1 (SD2)

Figure 6

- Björk Explains Decision To Pull <i>Vulnicura</i> From Spotify
- "As Punisher Joins <i>Daredevil</i> Season Two, Who Will the New Villain Be?"
- Talks on the Precepts and Buddhist Ethics

Figure 7

- Emma Watson to play Belle in Disney's <i>Beauty and the Beast</i>
- "<i>The Long Dark</i> Gets First Trailer, Steam Early Access"
- "Long-Lost F. Scott Fitzgerald Story Rediscovered and Published, 76 Years Later"

Figure 9

- Brit Marling-Zal Batmanglij Drama Series <i>The OA</i> Gets Picked Up By Netflix
- Sarah Silverman Will Star in HBO Pilot from <i>Secret Diary of a Call Girl</i> Creator
- Ethan Hawke to Star as Jazz Great Chet Baker in New Biopic
- "George R.R. Martin to Focus on Writing Next Book, World Rejoices"
- Hawkgirl Cast in <i>Arrow</i>/<i>Flash</i> Spinoff Series For The CW
- South Park: The Stick of Truth Review (Multi-Platform)
- Ava DuVernay Won't Direct <i>Black Panther</i> After All
- Freddy Adu Signs For Yet Another Club You Probably Don't Know
- "....."Listen to The Dead Weather's New Song, "....."Buzzkill(er)....."
- Emma Watson Set to Star Alongside Tom Hanks in Film Adaptation of Dave Eggers' <i>The Circle</i>

Figure 10

- Mothers influence on her young hippo
- Rambo 5 und Rocky Spin-Off - Sylvester Stallone gibt Updates
- Living in the Light with Ann Graham Lotz

- Brit Marling-Zal Batmanglij Drama Series <i>The OA</i> Gets Picked Up By Netflix
- Watch the First Episode of <i>Garfunkel and Oates</i>
- Emma Watson Set to Star Alongside Tom Hanks in Film Adaptation of Dave Eggers' <i>The Circle</i>
- Full body U-Zip main opening - Full body U-Zip main opening on front of bag for easy unloading when you get to camp
- Insights with Laura Powers
- Ethan Hawke to Star as Jazz Great Chet Baker in New Biopic
- Freddy Adu Signs For Yet Another Club You Probably Don't Know

Figure 12

MS COCO

- A painted sign of a blue bird in a tree in the woods.
- a small white toilet is in a tiny bathroom
- Two horses standing behind fence with grassy leaves
- A lot of wooden shelves filled with lots of clutter.
- A zebra standing on the grass holding its head near the ground.
- a skier and a snowboarder in front of a large house
- Large airplane flying below the clouds from underneath
- A snowboarder in a colorful jacket racing down a slope.
- There are police men standing and some are sitting on horses
- A large red vase sitting in front of a building.

Lexica Art

- Baby Yoda by Jeszika Le Vye, trending on artstation, hyperdetailed,
- the dnd monster pung_ concept art
- graffiti on a wall of anime gundam, speed painting, trending on artstation
- A beautiful concept art painting of a gloomy lake with a body floating slightly below the surface by
- three different views of a predator set for dota 2, concept art by senior character artist, trending on artstation, artstation hd, full body
- closeup of a pug with moss growing out of its face folds, macro photography, overgrown pug, high resolution photo, trending on artstation
- two beautiful idols standing face to face, trending art-

- station, pixiv, detailed anime art
- guild wars 2, cinematic battlefield, Hyperrealistic CGI Photorealistic cinematic volume lighting fanart on Art-Station full sun shine day concept art, digital art, high detail by Daniel Dociu
- minecraft style concept art of a blue portal to a fantasy world
- fantasy art of a bustling tavern in china, at night, by justin gerard, highly detailed digital art, trending on artstation

LAION

- kitchen ideas with islands simply home designs home design ideas 3
- 78 most terrific supple duvet sets with matching let curtains covers curtain duvets cover definition define snazzy aequorea plus thresh expensive quilts
- India begins next phase of COVID vaccination from Monday:
- Gary Hart Photography: Bright Angel Lightning, Grand Canyon
- Beautiful happy lady drinking glass of red wine.
- berries_lead
- 2018 Ford F-150
- Alpine Falls in Loyalsock State Forest
- Retro Classic Full Metal Tear Drop Clear Lens Aviator Glasses C936
- Well M1911 A1 Vollmetall Springer 6mm BB Two-Tone

Figure 13

- 33 Screenshots of Musicians in Videogames
- “communication, email, mail, message, online, open icon”
- Anzell Blue/Gray Area Rug by Andover Mills
- Vector Tyre Icons isolated on white background Illustration
- *Bloodborne* Video: Sony Explains the Game’s Procedurally Generated Dungeons
- 3D Black & White Skull King Design Luggage Covers 007
- “document,print,preview icon”
- 12x12x1 Pure Green AC Furnace Air Filters Qty 6 - Nordic Pure

C. Difference of embeddings between CLIP and OpenCLIP

In the CLIPText encoder used in Stable Diffusion v1.4, text embeddings are grouped into three clusters: \mathbf{v}^{sot} , \mathbf{v}_i^{pr} , and a merged cluster of \mathbf{v}^{eot} and $\mathbf{v}_i^{\text{pad}}$. This is because the tokenizer duplicates the $\langle \text{eot} \rangle$ for padding, making $\mathbf{v}_i^{\text{pad}}$ nearly identical to \mathbf{v}^{eot} after encoding. In contrast, the CLIPText encoder of Stable Diffusion v2.1 yields a separate cluster for \mathbf{v}^{eot} and $\mathbf{v}_i^{\text{pad}}$, indicating a semantic distinction between the two.

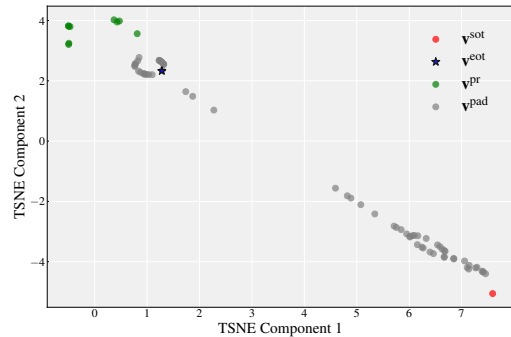
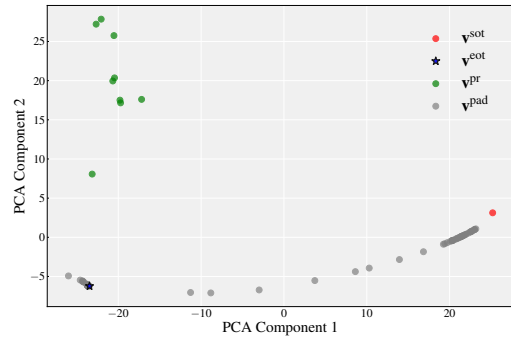
Stable Diffusion v1.4 employs the original CLIP text encoder, where the tokenizer uses $\langle \text{eot} \rangle$ as the default $\langle \text{pad} \rangle$. As a result, $\mathbf{v}_i^{\text{pad}}$ are near duplicates of \mathbf{v}^{eot} , which amplifies the semantic influence of \mathbf{v}^{eot} and contributes significantly to memorization.

Stable Diffusion v2.1 employs an OpenCLIP text encoder. Importantly, OpenCLIP maintains the same overall architecture and embedding mechanism as the original CLIP, with only minor differences in training scale and data. The critical difference relevant to our analysis is that the OpenCLIP tokenizer uses a neutral $!$ token for $\langle \text{pad} \rangle$ instead of duplicating $\langle \text{eot} \rangle$. As a result, $\mathbf{v}_i^{\text{pad}}$ are semantically distinct from the \mathbf{v}^{eot} , eliminating the structural duplication that contributed to memorization in Stable Diffusion v1.4.

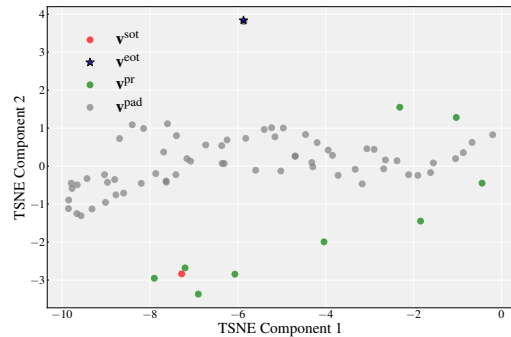
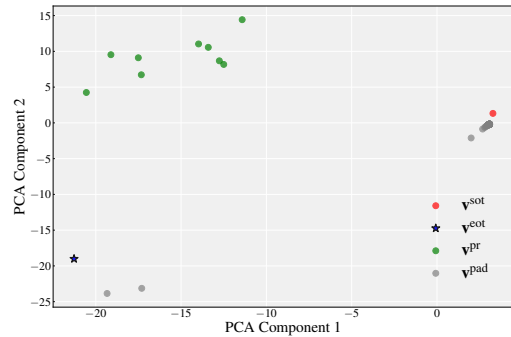
The PCA and t-SNE visualizations in Figure 14 further illustrate this difference.

D. Cross-Attention Maps by Embedding

Following Chen et al. [4], we visualize the cross-attention maps from the first two 64-pixel downsampling layers of the U-Net at the last denoising step. To ensure consistent brightness across all maps, we normalize the intensity by setting the minimum and maximum brightness levels equal across all cross-attention maps. As shown in Figure 15 and 16, not only does the attention map corresponding to \mathbf{v}^{eot} appear bright, but the adjacent $\mathbf{v}_i^{\text{pad}}$ also exhibit high attention scores.



(a) CLIP: PCA (top) and t-SNE (bottom)



(b) OpenCLIP: PCA (top) and t-SNE (bottom)

Figure 14. Text embeddings of “Living in the Light with Ann Graham Lotz” visualized using PCA and t-SNE for (a) CLIP and (b) OpenCLIP models. The two projection methods reveal distinct clustering behaviors.

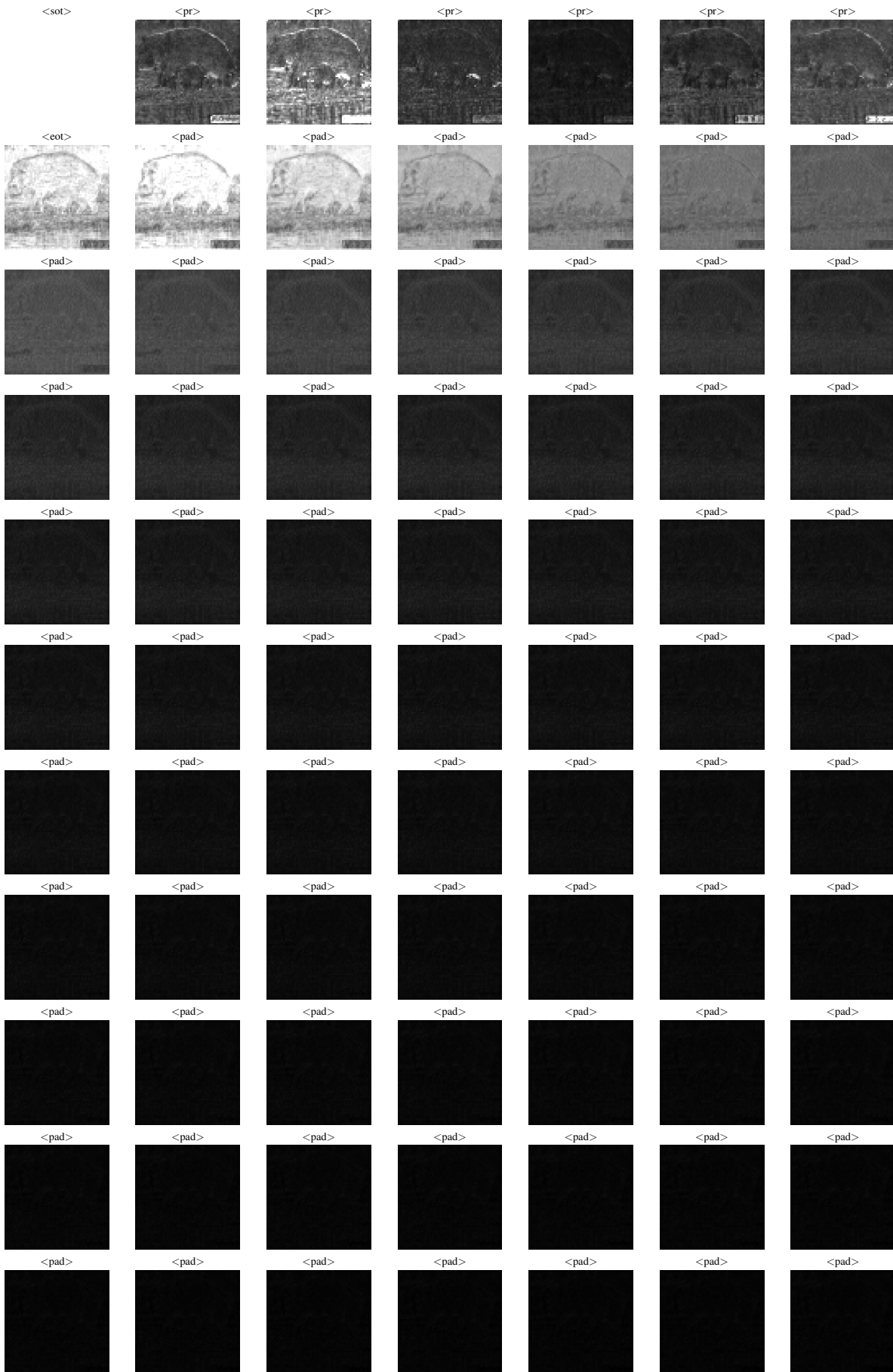


Figure 15. Cross Attention maps of “Mothers influence on her young hippo”

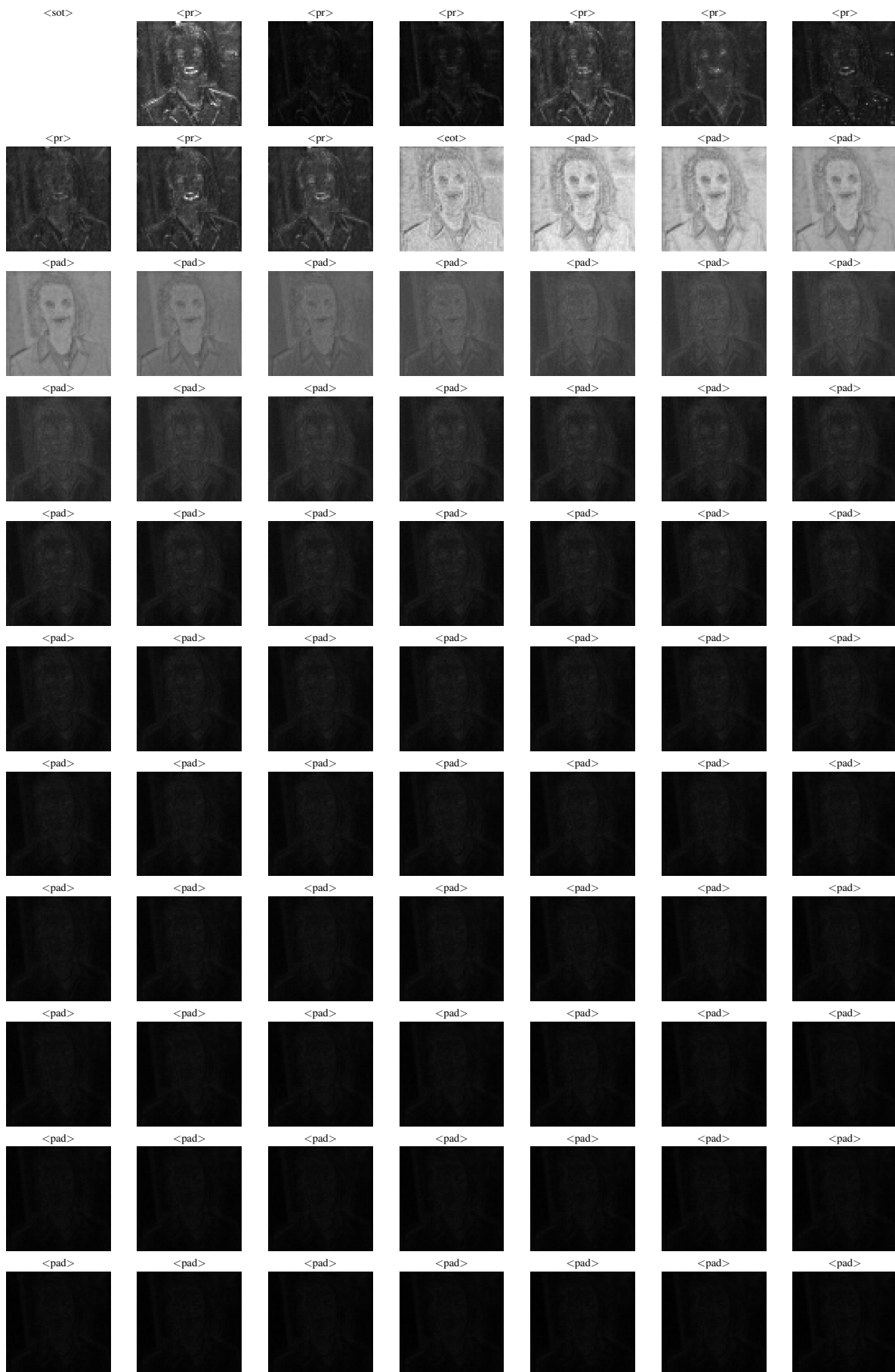


Figure 16. Cross Attention maps of "Living in the Light with Ann Graham"