053

Black Hole-Driven Identity Absorbing in Diffusion Models

Supplementary Material

In Appendix A, we delve into the Black Hole metaphor, 001 illustrating its connection to the BIA framework and ex-002 ploring identity distribution properties in high-dimensional 003 latent space, which underpin black hole region formation. 004 Appendix B empirically validates the boundary construc-005 006 tion for the black hole region and provides hyper-parameter configurations for optimal efficiency. Appendix C show-007 cases additional qualitative results, emphasizing the effec-008 tiveness of our method and some additional ablation study. 009 Finally, Appendix D outlines the complete algorithm, offer-010 ing a comprehensive overview of the BIA framework. 011

012 A. Black Hole Region Formation

013 This section provides a detailed explanation of the intuition
014 behind forming the black hole region in the latent space,
015 focusing on the efficient selection of identity-similar and
016 dissimilar points to ensure robust boundary formation.

017 A.1. Black Hole Metaphor

To form a black hole, you first need stars, as black holes 018 originate from the deaths of stars [7]. Similarly, in our prob-019 lem, we draw a parallel between the latent space of diffu-020 sion models and the universe. Just as stars populate the uni-021 verse, identity representations inhabit the latent space. To 022 create a black hole in this metaphorical universe, we des-023 024 ignate the dying identity, the one we aim to unlearn, as the black hole. This black hole absorbs nearby identity repre-025 sentations, gradually erasing its own identity features. 026

Taking inspiration from "Hubble discovers collection of 027 black holes in globular cluster" like NGC 6397 Fig. 1, 028 029 where stars are densely packed around a central black hole [3, 9], the latent space similarly exhibits clusters of iden-030 tity points that are naturally close to one another. In 031 this analogy, the surrounding identities act as neighboring 032 stars, which are pulled toward the black hole and ultimately 033 merged to remove the specified identity. 034

A.2. Identities Distribution in Latent Space

In forming the black hole region, we assume that identity-036 similar points are located in close proximity within the la-037 tent space. Our experiments on the CelebA-HQ dataset [2] 038 039 clearly illustrate this phenomenon, as mapping 20 identities with 10 images each demonstrates that similar identities 040 naturally cluster together within the latent space Fig. 2. This 041 supports the notion that identity representations are inher-042 ently structured, enabling us to define a precise black hole 043 region. By forming this region, we effectively pull in sur-044 045 rounding identities, facilitating the removal of the specified



Figure 1. Globular cluster NGC 6397, captured by NASA's Hubble Space Telescope, showcasing densely packed stars surrounding a central mass of black holes. (Image credit: NASA, ESA, T. Brown, S. Casertano, and J. Anderson, STScI) [3].



Figure 2. Latent space mapping between images and identities using t-SNE [8]. Each color represents a distinct identity with 10 images per identity from the CelebA-HQ dataset. The tight clustering of points highlights the proximity of similar identities, validating the Black Hole region's effectiveness in our BIA framework for identity unlearning.

identity. The validity of the relationship between CelebA-
HQ identities and its distribution in the latent space is also
examined by GUIDE [4], which confirm the natural close-
ness of identities in latent space. Much like astrophysicists
study the number and distribution of black holes in the uni-
verse [6], our approach models the latent space as a dynamic
environment, revealing its potential for unlearning identity.046
047

B. Boundary Validation:

We determine the boundary in the h-space by exploring the 054 neighborhood of the latent code h_r , representing the iden-055

069



Figure 3. Additional results on the CelebA-HQ dataset showcasing different identities. The first row represents the original identity, while the subsequent row illustrates the unlearned identity achieved through our BIA framework.

tity to be unlearned, similar to [5], to ensure a balanced 056 representation of identity-similar and dissimilar points for 057 training a robust SVM classifier. Through experimentation 058 059 with different combinations of n (number of latent points) and α (scaling factor for neighborhood exploration), we 060 identified the optimal range of n to be between 100-150 061 (noted as n in the main paper, with a typo in line 262) and 062 α between 0-20. These values effectively capture sufficient 063 identity-similar and non-similar points in appropriate pro-064 065 portions for classifier training. The resulting SVM achieves over 90% test accuracy, validating the identified hyperplane 066 d_{id} as the optimal boundary for identity disentanglement 067 within the latent space. 068

C. Additional Qualitative Results:

We present supplementary qualitative results using 8 im-
ages per identity, as shown in Fig. 3. These results reaffirm070that BIA effectively unlearns the target identity across both
the source and other images sharing the same identity while072preserving other attributes well.074

In Fig. 5, we demonstrate BIA's ability to maintain the 075 overall generative quality as outputs from the unlearned 076 model are visually consistent with those from the original 077 model. To ensure a fair comparison, the same seed was used 078 to generate images before and after model unlearning. The 079 preservation loss plays a crucial role in minimizing distribu-080 tion shifts, enabling BIA to effectively unlearn the specified 081 identity without compromising the quality of other latent 082 codes. 083

097

CVPR 2025 Submission #18418. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 4. Unlearned images from input images with increasing value of k in Random method.



Figure 5. Qualitative comparison showcasing the consistency of BIA with the original model. The unlearned model preserves the generation quality on randomly generated images, producing nearly identical outputs to the original model.

084 C.1. Ablation Study

In Section 4.3 of our main paper, we outlined the Ran-085 086 dom approach, which selects non-identity points to wrap the black hole region. Here, we investigate the effect of increas-087 ing the number of random points used for wrapping denoted 088 as k_r . As k_r increases, the Random approach introduces 089 090 more latent space noise, leading to a pronounced degradation in detail and structure. This progression is visualized in 091 Fig. 4, where the first column represents the original image, 092 and the column labeled Ours showcases the results of the 093 BIA framework. The subsequent four columns illustrate the 094 095 Random approach with increasing k_r values, demonstrating 096 a clear decline in visual quality and detail preservation.

D. Algorithm

The latent space of diffusion models exhibits a diverse iden-098 tity distribution, making it essential to define the black hole 099 region comprehensively to encompass all latent points as-100 sociated with a specified identity. To achieve this, we adopt 101 an efficient neighborhood sampling approach inspired by 102 GUIDE [5], but with key modifications for our framework. 103 The overall process of our BIA framework is explained 104 in Algorithm 1. As starting with h_r as the central refer-105 ence point, we iteratively sample neighboring latent codes 106 ${h_i}_{i=1}^n$ within its vicinity as we are only focusing on the 107 h_r to find the neighborhood points different form the [5], 108 which consider source and target latent both for finding 109 neighborhood points. These neighboring points are gener-110 ated by scaling the distance between h_r and randomly sam-111 pled latent codes hrandom, using a factor α drawn from a 112 uniform distribution $\mathcal{U}(0, \alpha_{\max})$, where α ranges from 0 to 113 20. This ensures sufficient local variation while maintaining 114 focus on the target identity. 115

The sampled latent codes are categorized as identity-116 similar or identity-dissimilar based on their cosine similar-117 ity to x_r , evaluated using arcFace feature embeddings [1]. 118 These labeled points are used to train an SVM, defining a 119 hyperplane that separates identity-similar points (inside the 120 black hole) from dissimilar points (outside). This process 121 ensures that the black hole region comprehensively repre-122 123 sents the target identity's latent distribution. To wrap the black hole region effectively, we incorporate $\mathcal{L}_{wrapper}$, which 124



a. Black hole region before unlearning b. Black hole region after unlearning

Figure 6. Illustration of Black hole region wrapping before and after unlearning.

minimizes the difference between the identity-similar and 125 dissimilar latent points, ensuring precise identity absorp-126 127 tion. Additionally, \mathcal{L}_{sem} is used to preserve non-identity attributes during optimization, ensuring high-quality outputs. 128 Together, these losses optimize the boundaries of the black 129 130 hole region while preserving the overall semantic integrity of the latent space. 131

Through this process, we fine-tune the diffusion model to 132 effectively unlearn the specified identity while maintaining 133 the generative quality and preserving non-identity attributes 134 as visualized in Fig. 6. 135

References 136

- 137 [1] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for 138 deep face recognition. In Proceedings of the IEEE/CVF 139 conference on computer vision and pattern recognition, 140 pages 4690-4699, 2019. 3 141
- [2] Tero Karras. Progressive growing of gans for im-142 proved quality, stability, and variation. arXiv preprint 143 144 arXiv:1710.10196, 2017. 1
- [3] Paul Morris. Hubble discovers collection of black holes 145 in globular cluster, 2021. 1 146
- [4] Juwon Seo, Sung-Hoon Lee, Tae-Young Lee, Seungjun 147 Moon, and Gyeong-Moon Park. Generative unlearning 148 149 for any identity: Supplementary material. 1
- [5] Juwon Seo, Sung-Hoon Lee, Tae-Young Lee, Seungjun 150 Moon, and Gyeong-Moon Park. Generative unlearning 151 for any identity. In Proceedings of the IEEE/CVF Con-152 ference on Computer Vision and Pattern Recognition, 153 pages 9151-9161, 2024. 2, 3 154
- [6] Alex Sicilia, Andrea Lapi, Lumen Boco, Mario Spera, 155 Ugo N Di Carlo, Michela Mapelli, Francesco Shankar, 156 David M Alexander, Alessandro Bressan, and Luigi 157 Danese. The black hole mass function across cosmic 158 times. i. stellar black holes and light seed distribution. 159 The Astrophysical Journal, 924(2):56, 2022. 1 160
- [7] Paul Sutter. How many black holes are there in the uni-161 162 verse?, 2021. 1

Algorithm 1: Black Hole-Driven Identity Absorption (BIA)

Input: Diffusion model ϵ_{θ} , unlearning identity image x_r .

Output: Updated model $\bar{\epsilon}_{\theta}$ with x_r id unlearned.

- 1 Initialization: Extract latent code h_r for x_r using DDIM inversion.
- 2 Step 1: Black Hole Region Formation;
- **3** for j = 1, ..., n do
- Sample $h_{\text{random}} \sim \mathcal{H}$; 4
- Compute $\Delta_j = \alpha \frac{h_{\text{random}} h_r}{\|h_{\text{random}} h_r\|}, \ \alpha \sim \mathcal{U}(0, \alpha_{\text{max}});$ Define $h_j = h_r + \Delta_j;$ 5
- 6
- 7 Compute $sim_{r,j}$ using arcFace ϕ : $sim_{r,j} = \phi(x_r) \cdot \phi(\text{Decode}(h_j));$
- Classify h_j as identity-similar or dissimilar 8 based on $\sin_{r,j} > \operatorname{th}_r$;
- 9 Train an SVM to compute hyperplane d_{id} separating identity-similar and dissimilar points, forming the black hole region $\mathcal{B}_{h} \subset \mathcal{H}$.
- 10 Step 2: Black Hole Wrapping and Identity Absorption;
- 11 Find k-nearest latent points $\{h_j\}_{j=1}^k$ (also as h_k for simplicity) around \mathcal{B}_h , ensuring no overlap with h_r ;
- 12 Calculate wrapper loss:

 $\mathcal{L}_{\text{wrapper}} = \lambda_{\text{id}} \mathcal{L}_{\text{id}} + \lambda_{L_2} \mathcal{L}_{L_2} + \lambda_{\text{perc}} \mathcal{L}_{\text{perc}}.$

for $h_{\mathcal{B}_h} \in \mathcal{B}_h$ do

Utilize d_{id} to walk for identity-edited latent 13 code:

$$h_{\mathcal{B}_{h}} = h_{\mathcal{B}_{h}} + \Delta_{h}$$

where Δ is a step size in [100, 120];

Compute semantic loss: 14

$$\mathcal{L}_{\text{sem}} = \left\| h_{\mathcal{B}_{h}} - \hat{h}_{\mathcal{B}_{h}} \right\|_{2}$$

15 Optimize:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{wrapper}} + \lambda_{\text{sem}} \mathcal{L}_{\text{sem}}$$

for T iterations to update ϵ_{θ} , yielding $\overline{\epsilon}_{\theta}$. 16 Inference Evaluate $\bar{\epsilon}_{\theta}$ prevents x_r generation while preserving other attributes.

- [8] Laurens Van der Maaten and Geoffrey Hinton. Visu-163 alizing data using t-sne. Journal of machine learning 164 research, 9(11), 2008. 1
- [9] Mike Wall. Scientists find clump of black holes inside the heart of globular cluster (video), 2021. 1

165 166

167