

1045 Appendix

1046 A CONFIDENCE MATCHING LOSS

1047 We explore three key loss functions: cross-entropy loss[4, 44], max-
 1048 margin loss[43] and Poincaré loss[37]. Inspired by [29], we repre-
 1049 sent the model’s unnormalized scores by taking the product of the
 1050 weights of the last layer and the penultimate layer activations . The
 1051 cross-entropy loss aims to minimize the negative log-likelihood
 1052 of the identity under the model parameters. The formulation is as
 1053 follows:

$$\mathcal{L}_{ce} = -\log \left(\frac{\exp(p^T w_c)}{\exp(p^T w_c) + \sum_{j=1, j \neq c} \exp(p^T w_j)} \right) \quad (6)$$

1054 Here, p denotes the activation of the penultimate layer for sample
 1055 x , and w_c represents the weight of the last layer for the c -th class
 1056 in the target model M .

1057 The maximum margin loss not only encourages maximizing
 1058 the confidence score for the target class c but also emphasizes the
 1059 separability of this class from others. We reformulate the maximum
 1060 margin loss as follows:

$$\mathcal{L}_{mm} = -(p^T w_c) + \max_{j \neq c} (p^T w_j) \quad (7)$$

1061 Here, p represents the activation of the penultimate layer for sample
 1062 x . The weight of the last layer for the c -th class in the target model
 1063 M is denoted as w_c . The term $p^T w_c$ represents the unnormalized
 1064 logit value for the c -th class.

1065 Poincaré loss is a hyperbolic space embedding loss function,
 1066 which can be rewritten as:

$$\mathcal{L}_p = \operatorname{arccosh} \left(1 + 2 \frac{\left\| \frac{p^T w_c}{\|p^T w_c\|_1} - y_c \right\|_2^2}{\left(1 - \left\| \frac{p^T w_c}{\|p^T w_c\|_1} \right\|_2^2 \right) \left(1 - \|y_c\|_2^2 \right)} \right) \quad (8)$$

1067 Poincaré loss measures the distance between two vectors u and
 1068 v in the hyperbolic space. $\|\cdot\|_2$ represents the Euclidean norm,
 1069 satisfying $\|u\|_2 < 1$ and $\|v\|_2 < 1$. Here, $u = \frac{p^T w_c}{\|p^T w_c\|_1}$, $\|\cdot\|_1$
 1070 denotes the absolute value norm, v is the one-hot encoded vector
 1071 for class c , denoted as y_c , and we replace 1 with 0.9999. Poincaré
 1072 loss belongs to the hyperbolic distance learning paradigm, which
 1073 enables measuring and comparing distances between vectors in a
 1074 larger embedding space.

1075 B EXPERIMENTAL SUPPLEMENT

1076 B.1 Datasets

1077 CelebA is a dataset of celebrity face attributes that contains 202,599
 1078 face pictures of 10,177 celebrity identities. For the training data
 1079 of the target models (Resnet18, Resnet152, and DenseNet169), we
 1080 selected 1000 identities with the most number of samples, resulting
 1081 in a total of 30,038 images. The FaceScrub dataset provides cropped
 1082 face images of 530 identities. However, on the dataset’s official web-
 1083 site, instead of actual images, they provide download links for the
 1084 dataset. All identities are used as target dataset. For Stanford.dogs,
 1085 this dataset is built on top of ImageNet, which is intended for non-
 1086 commercial research purposes only and provides 20,580 images of
 1087 120 dog breeds. For all datasets, the input images for the target

1088 model are resized to 224x224, while the input images for the evalua-
 1089 tion model(Inception-V3) are resized to 299x299. The CelebA[25]
 1090 dataset comprises 202,599 face pictures of 10,177 celebrity identi-
 1091 ties. In the case of training the target models (Resnet18, Resnet152,
 1092 and DenseNet169), we specifically selected 1,000 identities with the
 1093 highest number of samples, resulting in a total of 30,038 images.
 1094 Please note that the FaceScrub[27] dataset provides cropped face
 1095 images of 530 identities, but the dataset’s official website only offers
 1096 download links instead of actual images. All identities from Face-
 1097 Scrub were utilized as a target dataset for our research. Regarding
 1098 the Stanford.dogs[23] dataset, it is constructed on the foundation
 1099 of ImageNet[7], which is exclusively intended for non-commercial
 1100 research purposes. The ImageNet dataset provides 20,580 images
 1101 encompassing 120 dog breeds. In our experiments, the input im-
 1102 ages for the target models were uniformly resized to dimensions of
 1103 224x224 pixels, ensuring consistency across all datasets. However,
 1104 it is important to note that for the evaluation model (Inception-V3),
 1105 the input images were resized to a different size of 299x299 pixels.

1106 B.2 Publicly Available Image Prior

1107 We downloaded the code for StyleGAN2[19] from the official source
 1108 at stylegan2-ada. The AFHQ dataset consists of 16,130 high-resolution
 1109 images with a resolution of 512x512 pixels. We obtained the pre-
 1110 trained model weights for AFHQ.dogs512[5] by using the provided
 1111 link from the official code of StyleGAN2-ADA: AFHQ.dogs. It should
 1112 be noted that the FFHQ[21] dataset comprises 70,000 high-quality
 1113 face images. Compared to CelebA and FaceScrub, the image quality
 1114 in FFHQ is significantly higher. Subsequently, we downloaded the
 1115 pretrained model weights for FFHQ256 from ffhq256 and for Met-
 1116 faces from Metfaces. During the attack process, for the StyleGAN2
 1117 model pretrained on FFHQ256, we performed central cropping to
 1118 200 and then resized the images to 224x224 before inputting them
 1119 into the target model. For the generated images on AFHQ.dogs, we
 1120 applied central cropping to 400 and then resized them to 224x224.
 1121 In the comparative experiments, we followed the same cropping
 1122 and resizing approach for other MIAs method.

1123 B.3 Additional Experimental Results.

1124 **Visual Comparison of Attack Results on Facescrub.** We have
 1125 also visualized the attack results using the FaceScrub dataset, com-
 1126 paring the performance of different methods. From the visualized
 1127 results in figure 6, our approach demonstrates the ability to gener-
 1128 ate synthesis images that better reflect the distinctive features of
 1129 the target model’s private training data.

1130 Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

1131
 1132
 1133
 1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160

