

APPENDIX A - ADDITIONAL DETAILS AND ANALYSIS

HARD- AND SOFTWARE DETAILS

This experiment was conducted on a server equipped with 4 Tesla P100 GPUs and 256GB of memory. All models in the experiment were implemented using PyTorch. The training process further relied on CUDA 11.4, Python 3.8.10, and PyTorch 1.12.1 with Torchvision 0.13.1.

IMPLEMENTATION DETAILS

We employ StyleGAN2 (Karras et al. (2020)) with pre-trained weights from the FFHQ256 dataset for style transfer. We trained three target models, ResNet-18 (He et al. (2016)), ResNet-152 (He et al. (2016)), and DenseNet-169 Huang et al. (2017), on the CelebA dataset. They achieved accuracies of 86.38%, 87.35%, and 85.39% respectively on the test set. The selected initial points z follow a Gaussian distribution with a dimensionality of 512. We evaluate the attack using Inception-v3 (Szegedy et al. (2016)) and FaceNet (Schroff et al. (2015)) as evaluation models. We utilize Adam optimizer (Kingma & Ba (2014)) with an initial learning rate of 0.001 and $\beta=(0.9, 0.999)$. The joint loss function is defined as $M_{gard} = \alpha_1 \mathcal{L}_2 + \alpha_2 \text{Cosine} + \alpha_3 \mathcal{L}_1$, where $\alpha_1=0.5$, $\alpha_2=0.5$, and $\alpha_3=0.5$.

Real Label	Label Restoration	
	Inference results	Label Restoration Accuracy
[569]	[569]	100.0%
[25,759]	[25,759]	100.0%
[405,556,587]	[405,556,587]	100.0%
[405,566,587,532]	[405,566,587,914]	75.0%
[768,996,199,28]	[768,199,224,87]	50.0%
[768,996,199,28,367]	[768,199,367,224,87]	60.0%
[768,996,199,28,367,390]	[768,199,367,390,224,87]	66.7%
[768,996,199,28,367,390,783]	[768,199,367,390,224,87,539]	57.1%
[768,996,199,28,367,390,783,765]	[768,199,367,390,765,224,87,539]	62.5%

Table 3: Accuracy in inferring multiple labels from the CelebA dataset.

ATTACK PARAMETERS

In this section, we provide a detailed description of the parameter settings used in FGL experiments. To present the information more clearly, I have organized it into a table, as shown in Table 4.

Experiment Name	Target Id	Batch Size	Multi-Seed	Init-Point	Epoch	Vgrad Parameter
Ablation experiment	[354, 788, 280, 556, 568]	5	5	3000→4	70	1
Joint loss function	random[30]	30	6	5000→6	70	1
Different network architectures	[354, 788, 280, 556, 568]	5	5	3000→4	70	1
Gradient Regularization	[556, 28, 379, 672, 81, 652, 718, 848]	8	2	3000→2	200	$10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 0.1, 1$
Comparison with the state-of-the-art	[556, 28, 379, 672, 81, 652, 718, 848]	1	2	3000→2	70	1
Different Batch Size	random[10,20,30,40,50,60]	10,20,30,40,50,60	6	5000→6	70	1

Table 4: In the experiment, we configured various parameters.

LABEL INFERENCE

In our experiments, we observed a phenomenon different from previous research: label inference poses unique challenges in the context of GIAs. Firstly, (Zhao et al. (2020)) is the first work to introduce label inference in GIAs, demonstrating high accuracy in inferring labels for individual samples. However, this approach is not suitable for large batches of data. To address this issue, (Yin et al. (2021)) proposed a label inference technique that is effective for large batches and validated its performance on the ImageNet dataset.

However, our experimental focus was on the CelebA dataset, consisting of 224×224 -pixel images of faces. We directly applied the batch label inference method proposed by (Yin et al. (2021)), but found it unsuitable for CelebA. While it performed well with a batch size 3 or smaller, achieving 100% inference accuracy, its performance degraded when the batch size exceeded 3, leading to a decrease in inference accuracy. We also investigated the inference performance for different labels,

as shown in Table 3. This had implications for our experiments. Consequently, when evaluating the effectiveness of our method, we opted to use ground truth labels directly rather than employing the label inference method.

APPENDIX B - ADDITIONAL EXPERIMENTS

RELATIONSHIP BETWEEN GRADIENT REGULARIZATION AND IMAGE QUALITY.

During the attack on the CelebA dataset, we encountered the issue of generated image gradients being too small, which made gradient matching challenging. To address this problem, we introduced the technique of gradient regularization in our research, which enables better optimization by processing gradients. However, we also observed that gradient regularization has an impact on image quality. We conducted a preliminary quantitative analysis of the gradient normalization parameter V_{grad} , and the results are presented in Table 5. Additionally, we showcase the image performance in Figure 8.

V_{grad}	Image Reconstruction Metric			
	Top-1 \uparrow	Top-5 \uparrow	D_{inc} \downarrow	D_{face} \downarrow
10^{-5}	1.00	1.00	1.00	0.71
10^{-4}	1.00	1.00	0.66	0.72
10^{-3}	0.00	1.00	0.86	0.86
0.01	1.00	1.00	0.41	0.86
0.1	1.00	1.00	0.58	0.71
1	1.00	1.00	0.46	0.62

Table 5: The Impact of Different V_{grad} Parameters on Image Quality.



Figure 8: The Impact of Different N_{grad} Parameters on Visual Image Quality.

In the experiment of gradient regularization technique, we tried different values of the regularization parameter V_{grad} ranging from 10^{-5} to 1. Surprisingly, setting higher values for the regularization parameter did not negatively impact the image quality. On the contrary, higher parameter values were found to be more favorable for the optimization process, leading to improved image quality and higher similarity with real images. This is because both the gradient ΔW and the pseudo-gradient $\Delta W'$ undergo normalization, which preserved the matching precision without significant

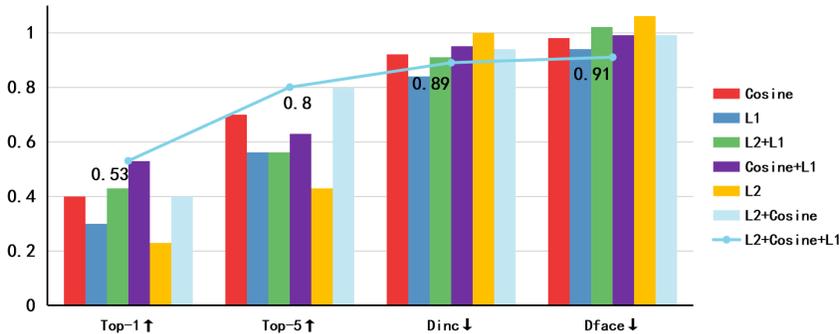
810 degradation. Consequently, in our experiments, we selected a regularization parameter $V_{grad} = 1$ for
 811 gradient normalization.
 812

813 M_{grad} JOINT LOSS FUNCTION.
 814

815 In the gradient matching task, Although we have demonstrated the effectiveness of the joint gradient
 816 function through ablation experiments, we are eager to gain a deeper understanding of its underlying
 817 mechanisms. We plan to further investigate how each individual regularization term in the joint
 818 gradient function contributes to the attack, whether a specific regularization term plays a crucial
 819 role independently, or if multiple regularization terms collaboratively generate the attack effect.
 820 Additionally, we intend to explore which part of the regularization terms significantly influences
 821 the attack results. we investigated the individual contributions of the components in the joint loss
 822 function $M_{grad} = \alpha_1 \mathcal{L}_2 + \alpha_2 \text{Cosine} + \alpha_3 \mathcal{L}_1$, as well as their combined effects. The results are
 823 shown in Table 6 and Figure 9.

M_{grad}	Image Reconstruction Metric			
	Top-1 \uparrow	Top-5 \uparrow	$D_{inc}\downarrow$	$D_{face}\downarrow$
\mathcal{L}_2	0.23	0.43	1.00	1.06
Cosine	0.40	0.70	0.92	0.98
\mathcal{L}_1	0.30	0.56	0.84	0.94
$\mathcal{L}_2 + \text{Cosine}$	0.40	0.80	0.94	0.99
$\mathcal{L}_2 + \mathcal{L}_1$	0.43	0.56	0.91	1.03
Cosine + \mathcal{L}_1	0.53	0.63	0.95	0.99
$\mathcal{L}_2 + \text{Cosine} + \mathcal{L}_1$	0.53	0.80	0.89	0.91

824
 825
 826
 827
 828
 829
 830
 831
 832
 833 Table 6: Comparing the Impact of Joint Gradient Matching Loss on Image Quality.
 834

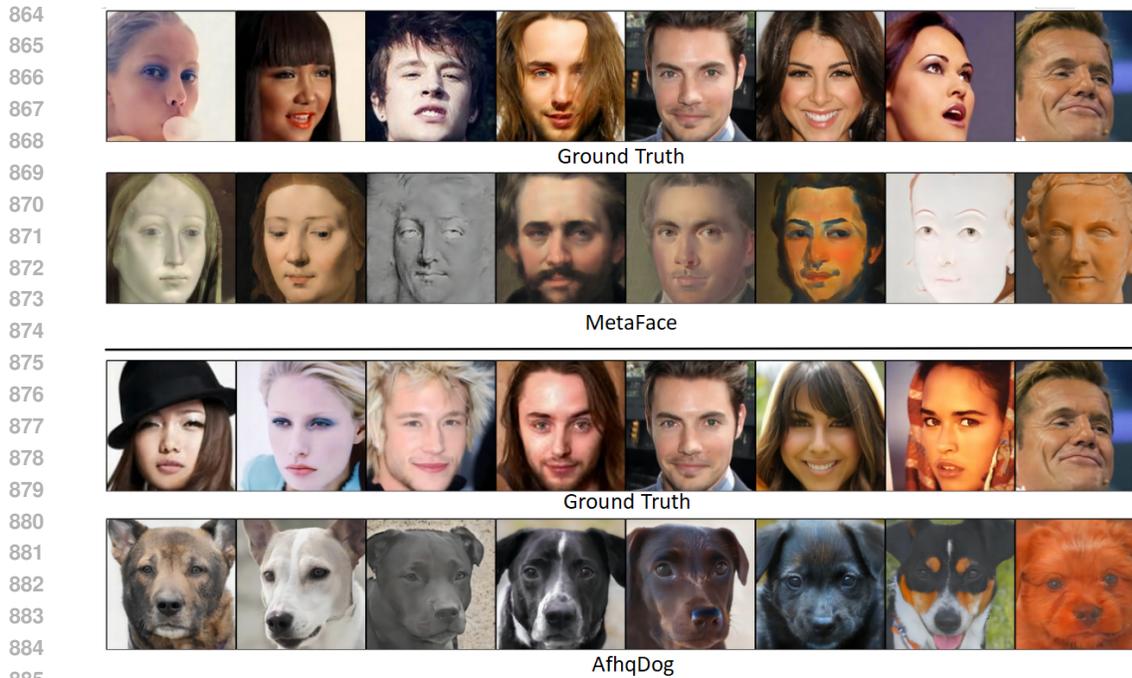


835
 836
 837
 838
 839
 840
 841
 842
 843
 844
 845
 846
 847 Figure 9: We conducted a study on the impact of different gradient matching loss functions on the
 848 overall attack effectiveness. In the ablation experiments, we used bar charts to illustrate the influence
 849 of different components, while for the complete joint function, we employed line charts to present
 850 the results.

851 In the experiments with the joint gradient loss function, we observed that each individual regulariza-
 852 tion term did not yield satisfactory results, with the cosine term slightly outperforming the others.
 853 However, when combining two regularization terms, the combination of cosine and \mathcal{L}_1 produced
 854 the best results. Surprisingly, the highest quality attack performance was achieved when all three
 855 regularization terms were combined simultaneously. Therefore, we believe that in the joint gradient
 856 loss function, the combined effect of multiple regularization terms plays a crucial role in achieving
 857 the optimal attack performance, rather than relying solely on individual effects or simple stacking
 858 of the terms.
 859

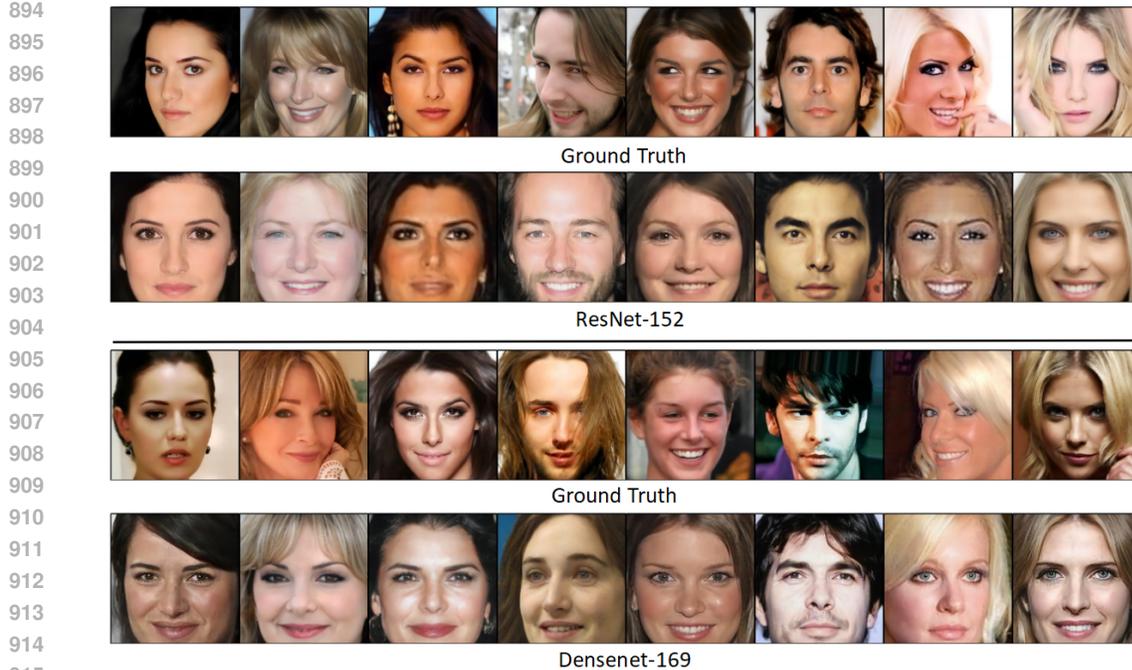
860 FGL ON DIFFERENT DATA DISTRIBUTIONS
 861

862 To validate the effectiveness of FGL in attacks under diverse data distributions, we not only con-
 863 ducted attacks on the Celeba dataset using the FFHQ dataset, but also performed attacks under ex-
 864 treme conditions using the significantly different MetFaces and AfhqDog datasets. The experimen-



886 Figure 10: We conducted experiments using datasets (MetFaces and AfhqDog) with entirely differ-
887 ent data distributions from CelebA.
888
889

890 tal results, as shown in Figure 10, demonstrate a noticeable resemblance between images generated
891 from both the MetFaces and AfhqDog datasets and the real images. This affirms that FGL is capable
892 of learning feature distributions similar to real data from datasets with distinct data distributions.
893



916 Figure 11: Experiments were conducted on deeper models, ResNet-152 and DenseNet-169.
917

918 FGL ON DIFFERENT NETWORK ARCHITECTURES

919
920 We not only verified the effectiveness of FGL on ResNet-18, but also conducted performance tests
921 on deeper models, ResNet-152 and DenseNet-169, to assess its practical applicability. As shown
922 in Figure 11, the experimental results are evident: both ResNet-152 and DenseNet-169 demonstrate
923 remarkable performance in generating images highly similar to real ones.

924
925 FGL ON DOG DATASET

926
927 In addition to performing well on the facial dataset, FGL was also validated on an animal dataset.
928 We trained ResNet-18 as the target model using the Stanford Dogs dataset and utilized Animal
929 Faces-HQ Dogs (AFHQ Dogs) as the image prior to train the GAN. The experimental results are
930 shown in Figure 14.

931 In our study, we observed an interesting phenomenon: humans are less sensitive to subtle differences
932 in objects or animals, but more sensitive to variations in facial images. This heightened sensitivity
933 towards facial images can be attributed to the frequent exposure to diverse facial representations in
934 daily life, making us more attuned to facial changes. In contrast, for small animals like dogs, minor
935 differences may not capture as much attention. This phenomenon is evident in the experimental
936 examples we provide.

937 FGL ON FEDAVG

938
939 Differing from FedSGD, FedAVG (McMahan et al. (2016)) performs multiple local updates before
940 sending the model weights w to the server. In our experiments, we conducted multiple trials of
941 attacks on FedAVG using FGL. In the configuration of FedAVG, we defined some hyperparameters
942 for local updates. Here, E (epoch) represents the local epoch, It (iteration) denotes the number of
943 local updates, and bs (batch size) indicates the batch size for each local update. We conducted
944 experiments with different parameters, and the results are shown in Figure 12 and Table 7 .

945

E/It/bs	Distance to Original Images			
	TOP-1	TOP-5	$D_{inc} \downarrow$	$D_{face} \downarrow$
E=1 It=8 bs=1	0.500	0.5	427	0.94
E=1 It=4 bs=2	0.375	0.625	398	1.00
E=2 It=8 bs=1	0.25	0.375	439	0.97
E=2 It=4 bs=2	0.5	0.5	449	1.14
E=3 It=8 bs=1	0.375	0.75	452	0.98
E=3 It=4 bs=2	0.25	0.25	502	1.10

953

954 Table 7: We conducted a series of FGL attack experiments with various parameters on FedAVG.

955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

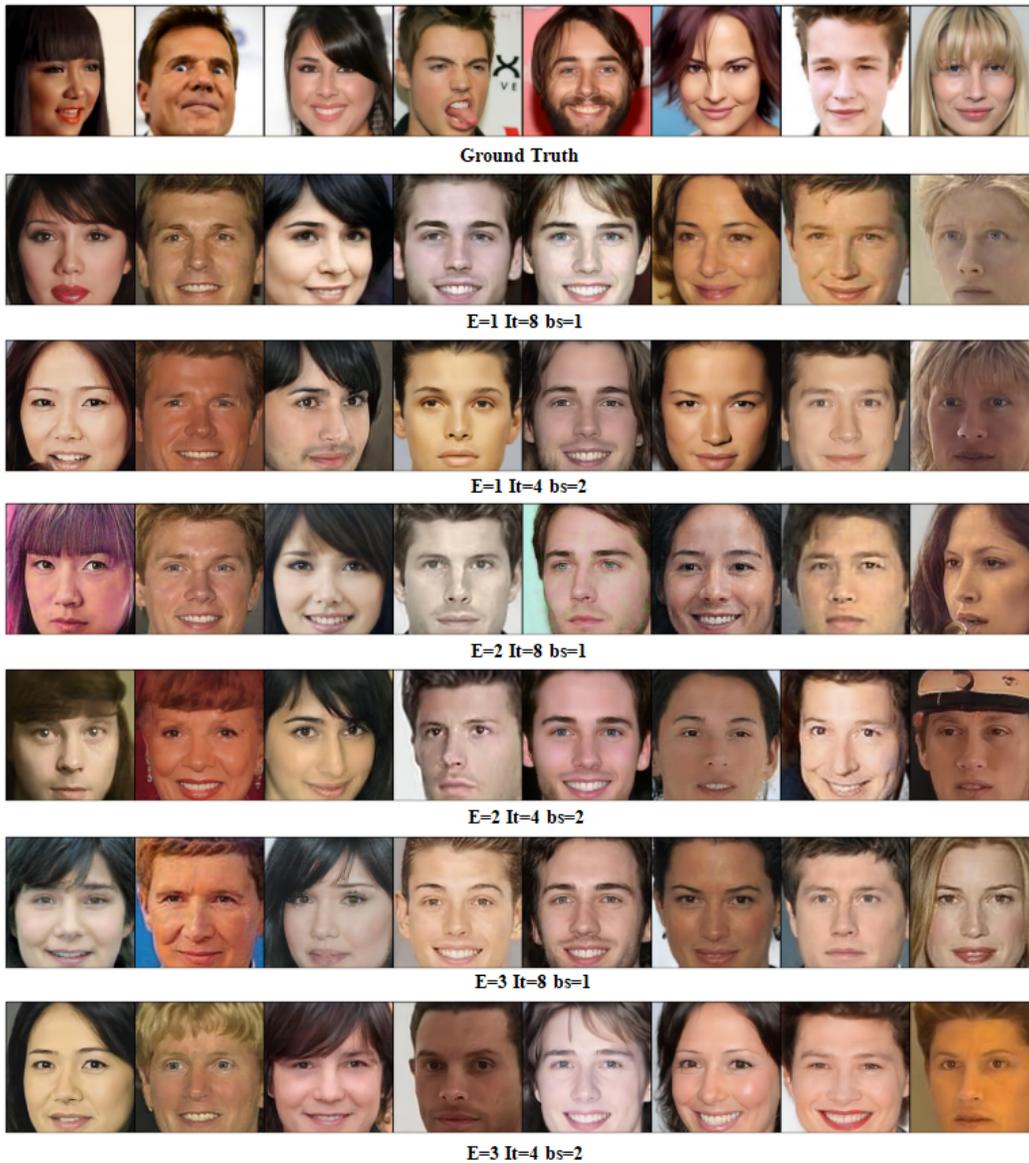


Figure 12: Here are the qualitative experimental results of FGL on FedAVG with different parameter settings.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

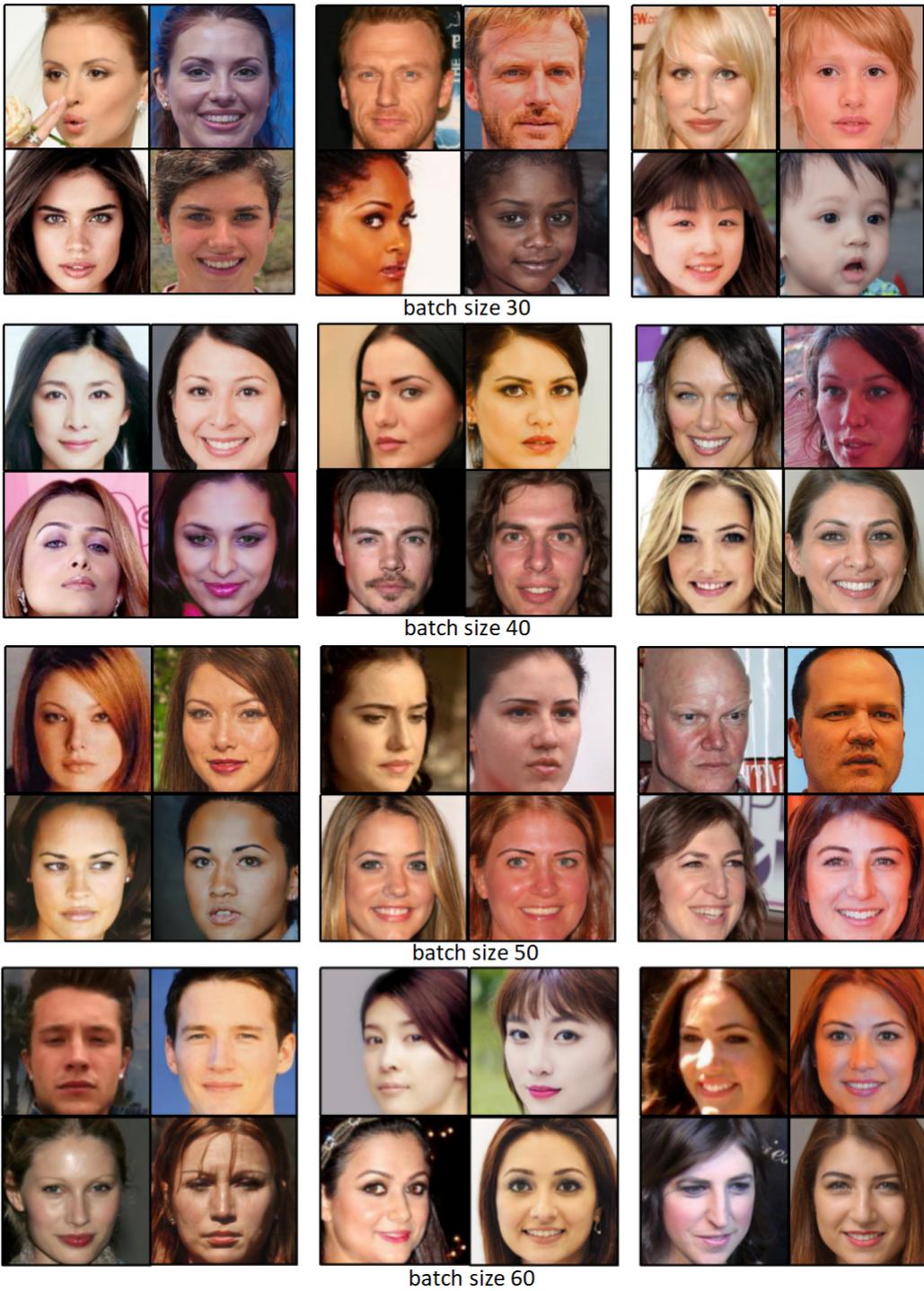


Figure 13: Provide more examples of large-batch attacks, with real images on the left and synthesized images on the right.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

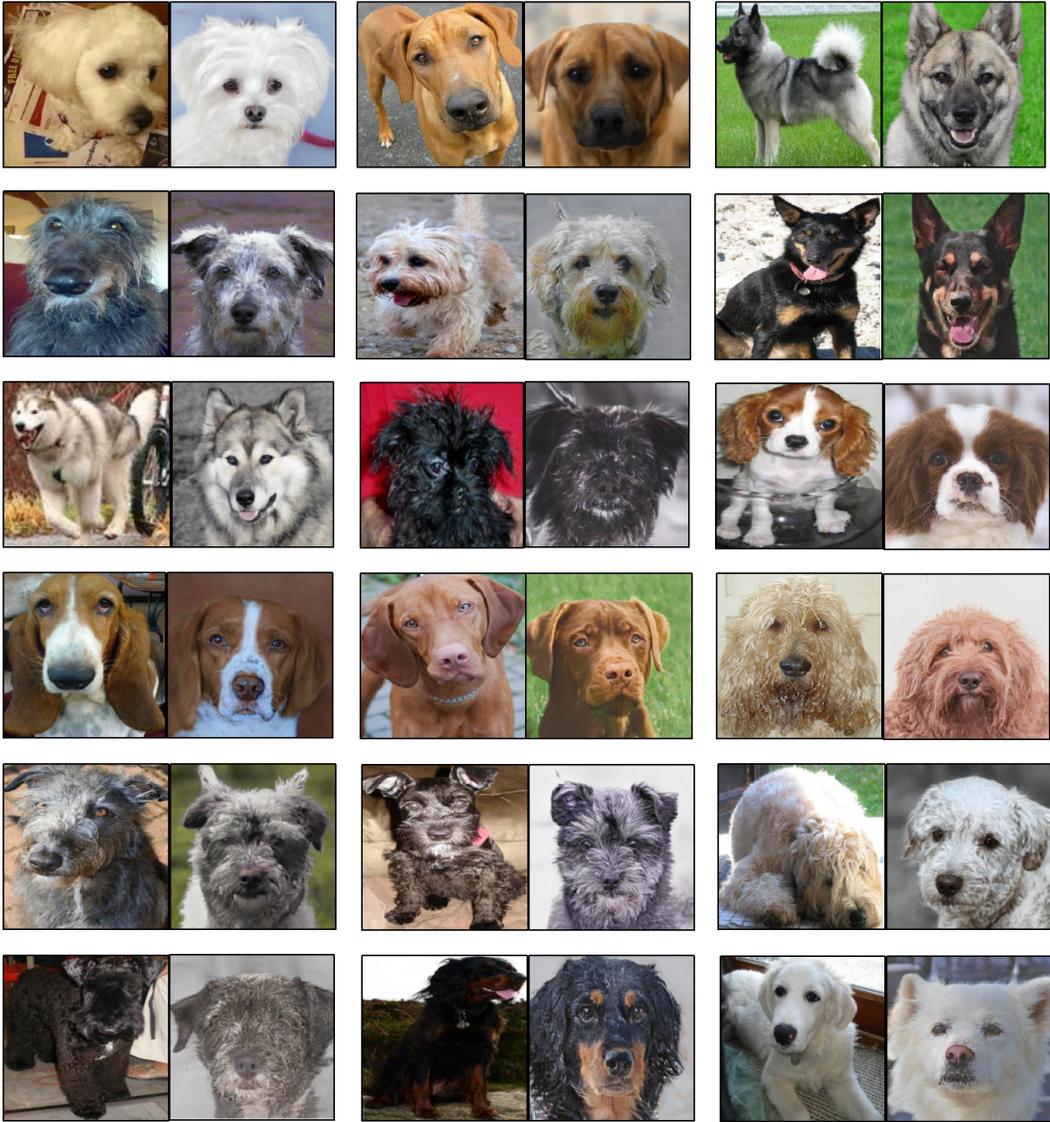


Figure 14: FGL conducted attacks on the Stanford Dogs dataset, with real images on the left and synthesized images on the right.