1 We thank the recognition from reviewers on the value of research problem, novelty of the proposed method and our
2 results. We address major raised concerns below. One missing reference (R1) will be added in the revised version.

3 **R1: fine-tuned or trained from scratch.** The work of [42] has shown that given enough data, training from scratch
4 or fine-tuning from a pretrained model could both learn a good generative model for the target domain, while the
5 fine-tuning simply gets faster convergence. As stated on L2-4, for the few-shot scenario, while it is unlikely to train the
6 model from scratch with a few data, we design our pipeline in the fine-tuning fashion. This is why we identify the trend
7 from the fine-tuned model because there are more correspondences on weights changes under the same way of learning.

8 **R1: FID on other shots.** We additionally evaluate our method against existing approaches on other shots in the
9 following table. The performance of NST will not be improved obviously because the style transfer method cannot
10 capture the target style well (L200-202). As stated on L244-246, the BSA's performance drops when increasing the
11 number of shots. The MineGAN and our work behave similarly as both are distribution learning based methods. The
12 more data, the better performance and the closer these two methods are.

Table 1: Quantitative comparisons between different few-shot generation methods (FID↓).

| Number of shots | NST [6] | BSA [28] | MineGAN [41] | Ours |
|---|---|---|---|---|
| 1 | $212.23 \pm 9.77$ | $102.34 \pm 5.70$ | $102.57 \pm 4.76$ | $84.36 \pm 3.91$ |
| 10 | $204.16 \pm 9.28$ | $105.56 \pm 5.79$ | $86.44 \pm 4.38$ | $74.87 \pm 3.75$ |
| 100 | $199.52 \pm 9.02$ | $110.24 \pm 5.87$ | $76.23 \pm 4.05$ | $67.55 \pm 3.48$ |
| 1,000 | $196.43 \pm 8.83$ | $119.31 \pm 5.96$ | $69.20 \pm 3.59$ | $62.40 \pm 3.12$ |
| 10,000 | $194.88 \pm 8.71$ | $131.20 \pm 6.04$ | $58.69 \pm 3.14$ | $55.74 \pm 2.88$ |

13 **R2: clarification on FI.** The Fisher Information is computed for each pretrained parameter $\theta_{S,i}$ in the **source** model,
14 and used as an importance weight to regularize the changes of each parameter $\theta_i$ during the adaptation on target domain.
15 The $i$ is the index of each parameter in the model. Note that the FI is computed for each single parameter instead of
16 each layer. What Figure 2(right) shows is the **average** FI of all parameters at each layer for easier visualization.

17 **R2: effect of each layer.** Thanks for the suggestion to ablate the importance of each layer and we will include it in the
18 revised draft. We mainly focus on the individual parameter because the FI is computed for each parameter.

19 **R2: source domain in Table 3.** As stated on L261, we select the FFHQ face dataset as the source domain.

20 **R3: other datasets.** In addition to results ($256 \times 256$) on face datasets, we also presented higher-resolution ($512 \times 512$,
21 $1024 \times 1024$) face generation in Figure 19-20 of the supplementary material and the landscape generation in Figure 4 of
22 the paper. We think our approach can be also used for other datasets with more structures details based on two premises:
23 i) a decent pretrained model on the source domain and ii) some level of similarity between the source and target domain.

24 **R3: more baselines on $F_i$.** The baseline of removing $F_i$ (i.e., without EWC used) is shown in Figure 3 of the paper
25 and discussed on L161-173. We also add another baseline by giving some fixed weights (i.e., the average FI in Figure 2
26 (right)) for all parameters in each layer. We observe that it will not lead to obvious over-fitting and the performance is
27 slightly worse (FID: 77.16 vs. Ours 74.87 for 10-shot). Comparing with treating all parameters at each layer equally
28 with the same weight, ours by regularizing each single parameter with its own FI still works better.

29 **R4: limited contribution.** Technically, our work on regularizing the model's own parameters is to **avoid** introducing
30 additional new parameters as did in existing approaches [28,41] which involves many tedious manual designs (e.g., the
31 number of parameters, the position to embed those parameters). The proposed method is simple and effective, and
32 may shed light on more future understandings of the learned parameters. Practically, we proposed a solution to the
33 challenging few-shot generation problem (e.g., 10 examples) as generative models often struggle in the low-data regime.

34 **R4: later layers are important.** Thanks for pointing out this improper statement. We shall not conclude parameters in
35 the later layer are all important based on the **average** FI. It should be some important parameters with much larger FI
36 that results in the higher average FI. We will redesign Figure 2 with the standard deviation for better visualization.

37 **R4: analysis of the regularization weight.** We assume R4 is referring to $\lambda$ in Eq. (3) as this is the only parameter we
38 empirically set the value for. We show its effect on the performance (FID) in the following table (10-shot). A large $\lambda$
39 tends preserve more style of the source and a small $\lambda$ may result in the some level of over-fitting. In addition, we find
40 that (i) if the source and target are more similar, select a larger $\lambda$, and (ii) if more target data is given, select a smaller $\lambda$.

| $\lambda$ | $5 \times 10^6$ | $5 \times 10^7$ | $5 \times 10^8$ | $5 \times 10^9$ | $5 \times 10^{10}$ |
|---|---|---|---|---|---|
| FID ↓ | $78.21 \pm 3.80$ | $75.62 \pm 3.74$ | $\mathbf{74.87 \pm 3.75}$ | $77.79 \pm 3.78$ | $80.21 \pm 3.82$ |