

1 Thanks all the reviewers for acknowledging our contributions and their valuable comments.

2 **To Reviewer #1 Q1:** More modalities? **R1:** Thanks. We add experiments on the NJU-ID (different resolutions) [1] and
3 the Multi-PIE (different poses) [2] datasets. Details of these datasets are introduced as follows.

- 4 • The NJU-ID dataset consists of 256 identities with one ID card image (102×126 resolution) and one camer-
5 a image (640×480 resolution) per identity. Considering the few number of images in the NJU-ID dataset, we
6 use our collected ID-Photo dataset (1000 identities) as the training set and the NJU-ID dataset as the testing set.
- 7 • The Multi-PIE dataset contains 337 persons with different poses. We use profiles ($\pm 75^\circ$, $\pm 90^\circ$) and frontal faces as
8 different modalities. 200 persons are used as the training set and the rest 137 persons are the testing set.

9 The examples of dual generation are shown in Fig. 1. For the recognition performance, on the NJU-ID dataset, we
10 improve Rank-1 by 5.5% (DVG 96.8% - Baseline 91.3%) and VR@FAR=1% by 6.2% (DVG 96.7% - Baseline 90.5%)
11 over the baseline LightCNN-29. On the Multi-PIE dataset, the Rank-1 of $\pm 90^\circ$ and $\pm 75^\circ$ is increased by 18.5% (DVG
12 83.9% - Baseline 65.4%) and 4.3% (DVG 97.3% - Baseline 93.0%), respectively. All experiments demonstrate the
13 effectiveness of our method in other modalities.

14 **Q2:** More ablations? **R2:** For the generation model, the ablations of $\mathcal{L}_{\text{dist}}$, \mathcal{L}_{ip} and \mathcal{L}_{div} have
15 been reported in Table-1. We add the ablation of \mathcal{L}_{adv} in Eq. (10). That is, on the CASIA
16 NIR-VIS 2.0 dataset, the Rank-1 decreases 0.5% if \mathcal{L}_{adv} is not used. For the recognition
17 model, the effect of $\mathcal{L}_{\text{pair}}$ in Eq. (13) can be found in Table-2 ('+DVG' means using $\mathcal{L}_{\text{pair}}$).
18 All ablations reveal that each component of our method is useful. Especially for \mathcal{L}_{ip} , $\mathcal{L}_{\text{dist}}$ and
19 $\mathcal{L}_{\text{pair}}$, the Rank-1 decreases 5.5%, 4.9% and 2.1% respectively on the ablations. Moreover, our method is not sensitive
20 to the trade-off parameters in a large range. Please see Reviewer-2' R2 for details.



Figure 1: The examples of dual generation on the ID-Photo and the Multi-PIE datasets.

21 **To Reviewer #2 Q1:** The relationship between DVG and PIM? Some related works? **R1:** Thanks. The noise in PIM is
22 to help recover invisible details. The generated 'many' faces are required to be consistent with one ground truth. Hence,
23 PIM is still a conditional image-to-image translation method. As mentioned in the introduction, it faces diversity and
24 uniqueness limitations. Differently, our method belongs to unconditional generation. That is, we generate diverse new
25 paired faces from noise, which alleviates the above two limitations. We will cite these related works in our paper.

26 **Q2:** How to assign hyper-parameters? A sensitivity analysis? **R2:** The hyper-parameters are
27 set by balancing the magnitude of each loss function. Fig. 2 presents the sensitivity studies of
28 λ_1 , λ_2 and λ_3 in Eq. (10). For α_1 in Eq. (13), when α_1 is set to 0.0025, 0.005, 0.01, 0.02 and
29 0.04, the Rank-1 is 98.9%, 99.1%, 99.2%, 99.2% and 98.8%, respectively. We can observe
30 that our method is not sensitive to these hyper-parameters in a large range. For example, the
31 Rank-1 only decreases 0.3% when λ_1 changes from 0.1 to 0.4.

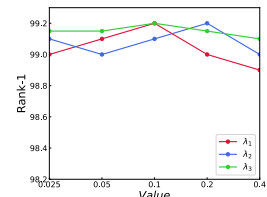


Figure 2: The sensitivity studies of trade-off parameters on the CASIA NIR-VIS 2.0 dataset. The backbone is LightCNN-9.

32 **Q3:** Complexity? **R3:** Thanks. Our method is computationally efficient. For instance, when
33 using one Titan XP, training the generation model on the CASIA NIR-VIS 2.0 dataset spends
34 3 hours. Meanwhile, in the inference stage, generating a pair of heterogeneous faces only
35 needs 3.2 ms. Furthermore, training the HFR network spends 1 hour.

36 **Q4:** Releasing the generated faces and the writing suggestions. **R4:** Thanks. We will release our codes and the
37 generated data. The writing of our paper has been carefully revised according to your advice.

38 **To Reviewer #3 Q1:** How to tune these three parameters in Eq. (10)? **R1:** The trade-off parameters in Eq. (10) are
39 tuned by balancing the magnitude of each loss function. In addition, the sensitivity studies of the trade-off parameters
40 λ_1 , λ_2 and λ_3 in Eq. (10) are shown in Fig. 2. We can observe that our method is not sensitive to these trade-off
41 parameters in a large range. For instance, when λ_1 changes from 0.1 to 0.4, the Rank-1 only decreases 0.3%.

42 **Q2:** It is suggested to report the time cost. **R2:** Thanks for your advice. Our proposed framework is computationally
43 efficient. For example, when using one Titan XP, training the generation model on the CASIA NIR-VIS 2.0 dataset only
44 needs 3 hours. Meanwhile, generating a pair of heterogeneous faces in the inference stage needs 3.2 ms. Moreover,
45 training the HFR network needs 1 hour.

46 **Q3:** Apply to other heterogeneous recognition problems? **R3:** We add experiments on other two datasets, including the
47 NJU-ID (different resolutions) [1] and the Multi-PIE (different poses) [2] datasets. The results show that our method
48 can be effectively applied to more modalities. Please see Reviewer-1' R1 for experimental details. Due to the limited
49 time, we will explore other heterogeneous recognition tasks in the future work.

50 **References:** [1] Huo et al. Heterogeneous face recognition by margin-based cross-modality metric learning. IEEE
51 Transactions on Cybernetics 2018. [2] Gross et al. Multi-PIE. Image and Vision Computing 2010.