# Improving the Transferability of Adversarial Attacks on Face Recognition with Diverse Parameters Augmentation

Fengfan Zhou[1], Bangjie Yin[2], Hefei Ling[1]*, Qianyu Zhou[3], Wenxuan Wang[4]

[1]School of Computer Science and Technology, Huazhong University of Science and Technology;

[2]Shanghai Shizhuang Information Technology Co., Ltd;

[3] Department of Computer Science and Engineering, Shanghai Jiao Tong University;
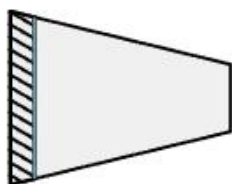
[4] School of Computer Science, Northwestern Polytechnical University.

[1]{ffzhou, lhefei}@hust.edu.cn, [2]jamesyin10@gmail.com,
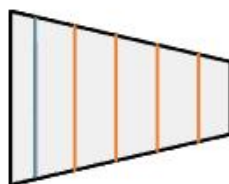
[3]zhouqianyu@sjtu.edu.cn, [4]wxwang@nwpu.edu.cn
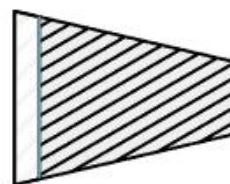
# Motivation

Pre-trained Parameters      Randomly Initialized Parameters
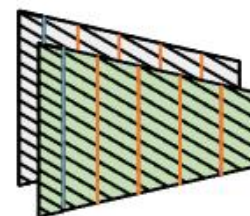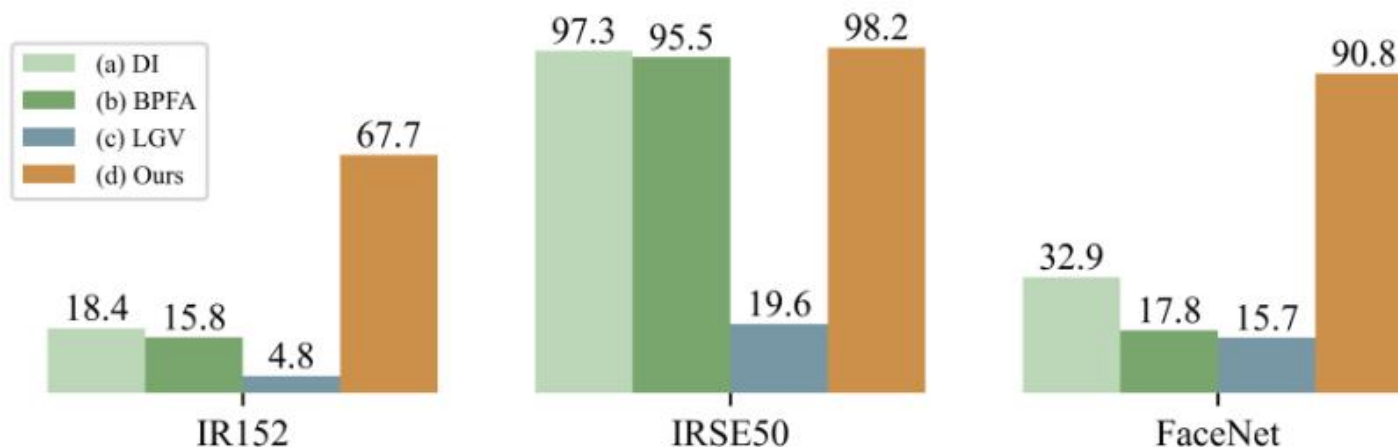
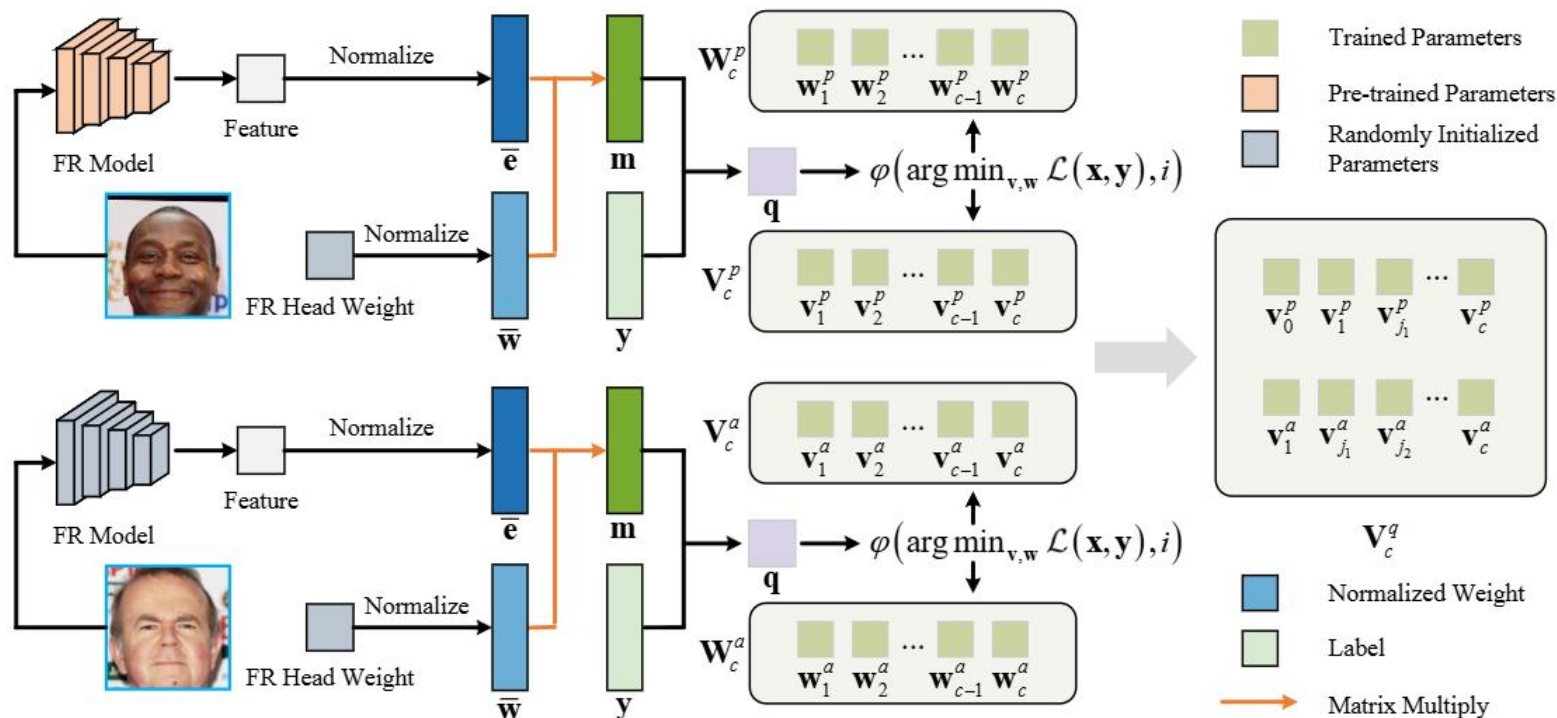(a) Input-based Augmentation

(b) Feature-based Augmentation

(c) Parameter-based Augmentation

(d) Ours

Legend:
- (a) DI
- (b) BPFA
- (c) LGV
- (d) Ours

IR152: 18.4, 15.8, 4.8, 67.7
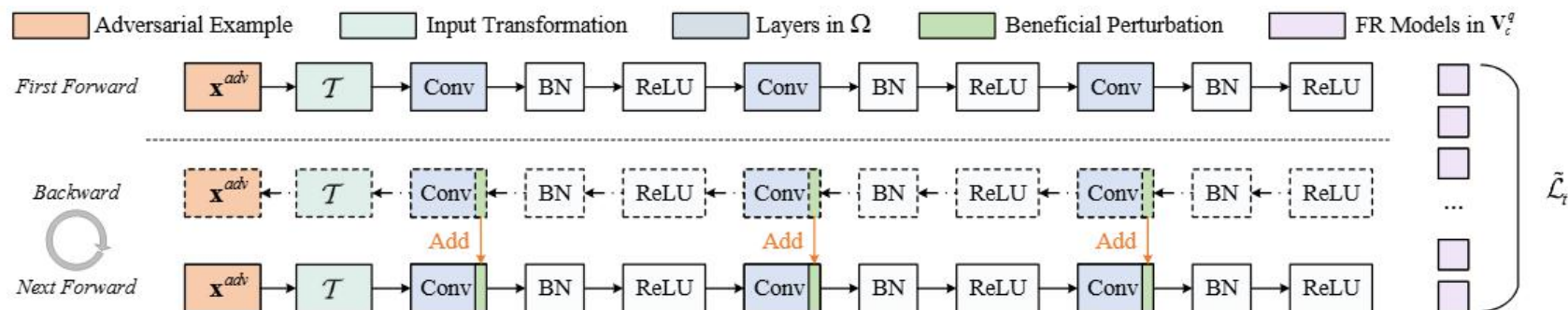
IRSE50: 97.3, 95.5, 19.6, 98.2

FaceNet: 32.9, 17.8, 15.7, 90.8

# Methodology

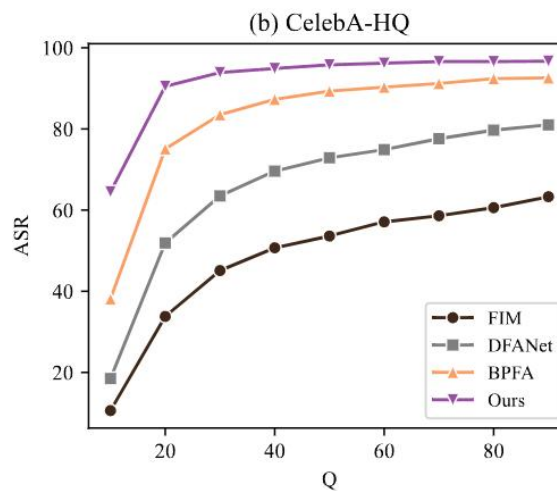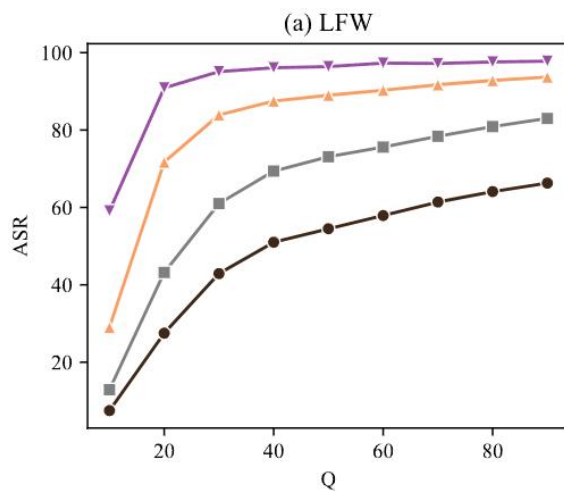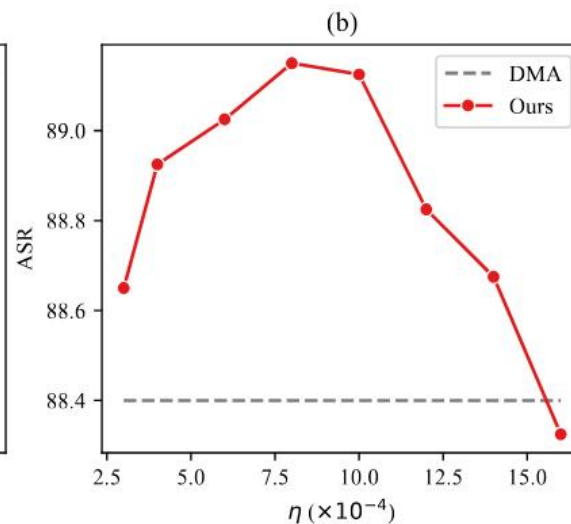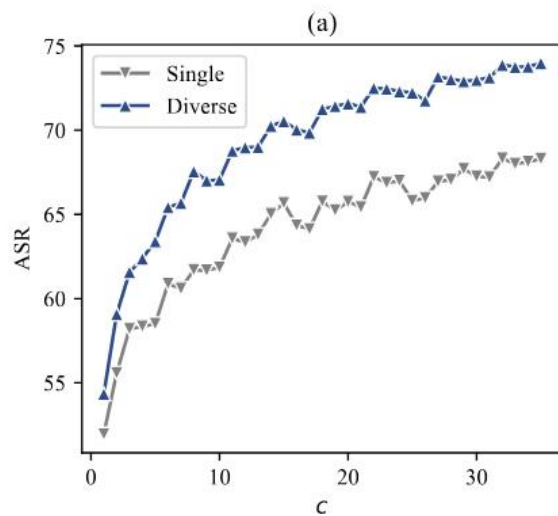## Diverse Parameters Optimization



We enhance the diversity of the surrogate model parameters by integrating both pre-trained and random initializations. The method yields a diverse set of surrogate model parameters, which enhances the parameter diversity of the surrogate FR models and consequently improves transferability of the crafted adversarial examples.

# Methodology

## Hard Model Aggregation



After acquiring a surrogate model set with diverse parameters, we introduce beneficial perturbations with the optimization direction opposite to that of adversarial perturbations onto the feature maps of these diversified surrogate models, transforming them into hard models and aggregate the hard models to increase the transferability.

# Experiments

# Thank you