# INTRODUCTION TO DEEPFAKE

Prepared by: Jet Kan, Yu Tianze
Advised by: Fang Cheng Fang
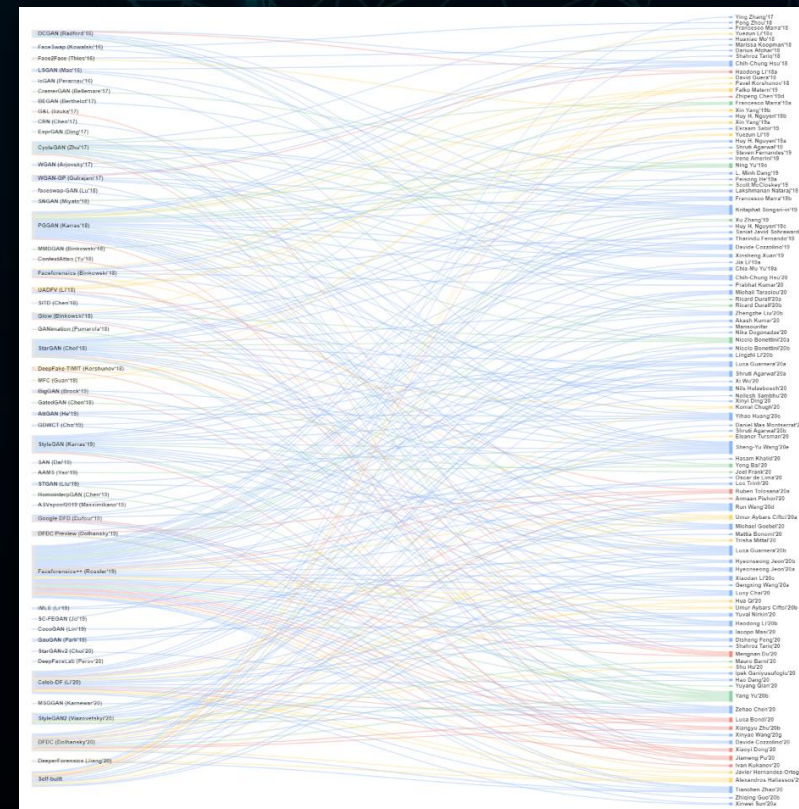
# THE DEEPFAKE BATTLEGROUND

Since its inception in 2016, rapid development of DeepFake in both generation and detection has formed the relationship of battleground, pushing the improvements of each other and inspiring new directions.
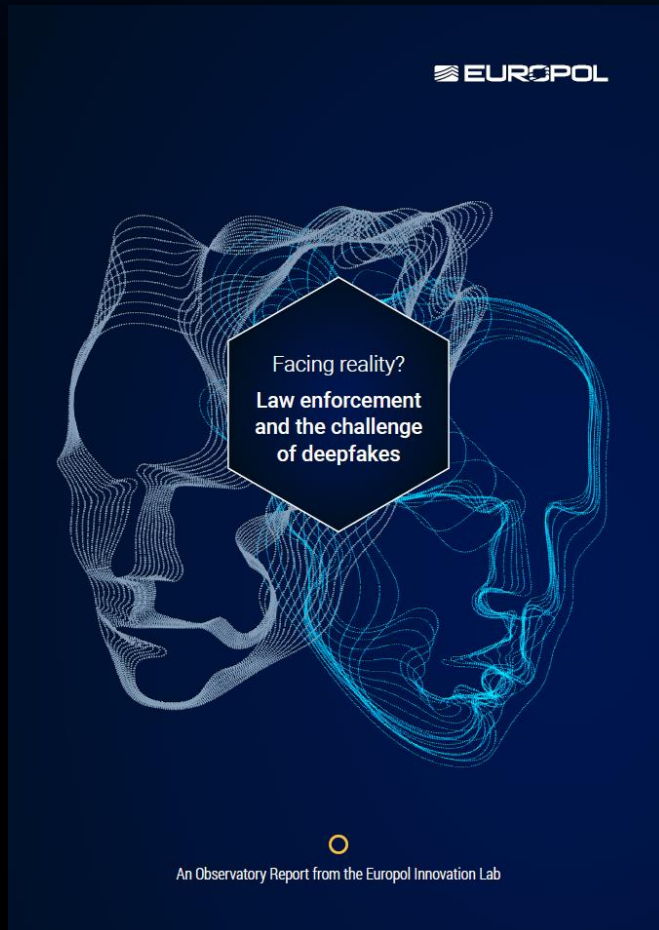


DF survey 2022



Interactions between generation and detection methods
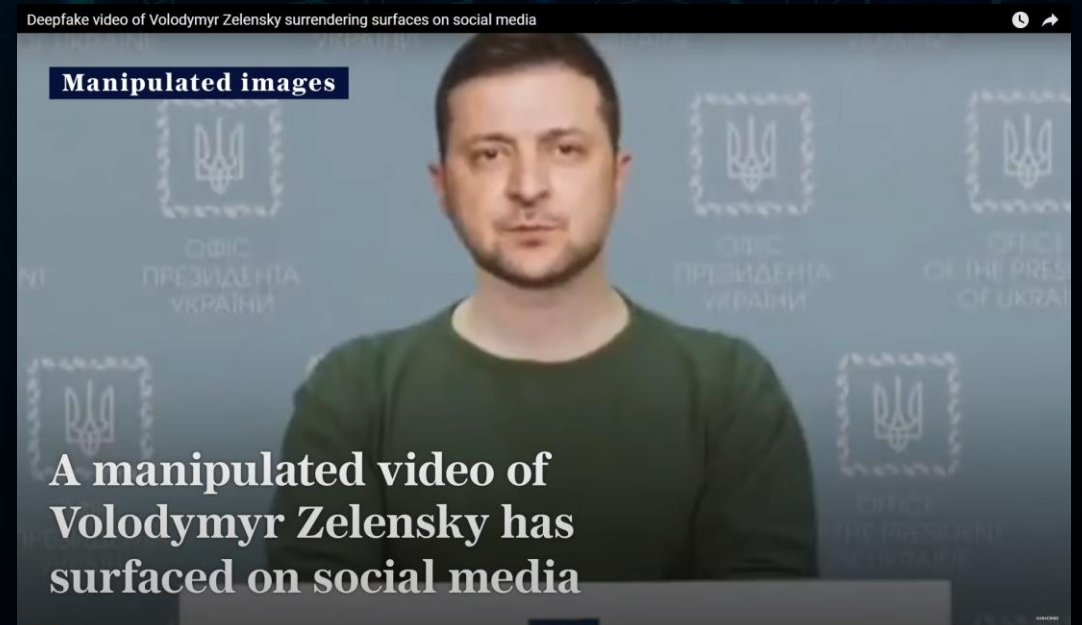
# EUROPOL INNOVATION LAB

- Europol: mandated by the EU to support the law enforcement in innovation

- Their first report published in 28 April 2022 explores the topic of DeepFake

- Summary:
  - Layman's introduction to DeepFake and Machine Learning
  - Impacts of DeepFake in law, financial and security
  - Current policies to combat fake videos

*"Experts estimate that as much as 90% of online content may be synthetically generated by 2026."*
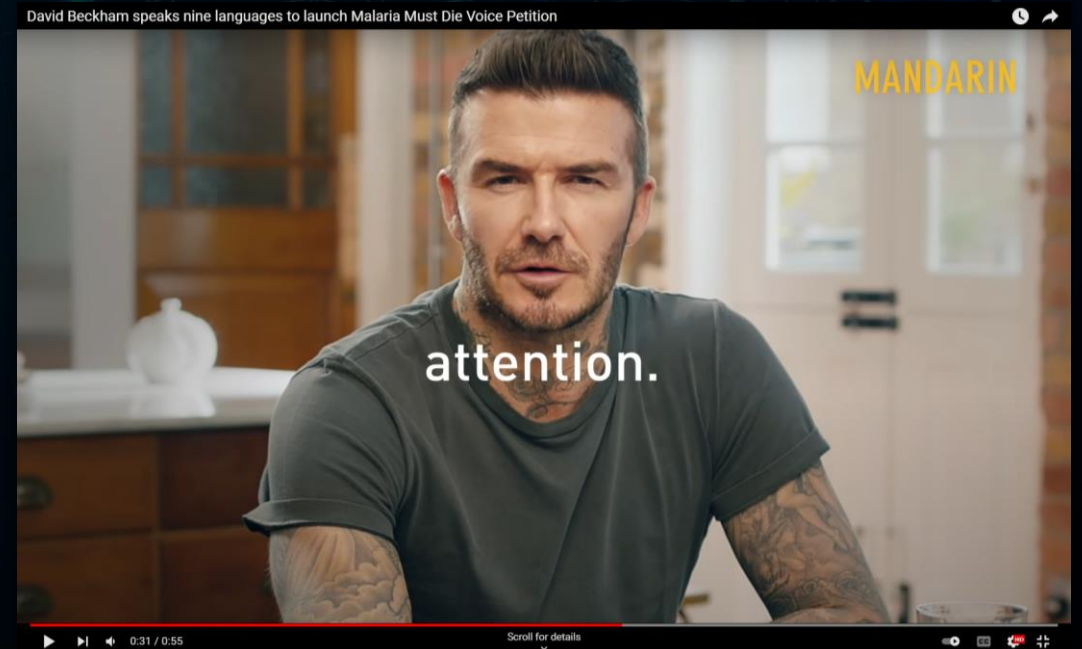
*- Europol (2022)*

# MALICIOUS USES OF DEEPFAKE

- Non-consensual pornography

- Document fraud

- Falsifying evidence for criminal justice investigations

- Distributing disinformation and manipulating public opinion

- Supporting the narratives of extremist or terrorist groups



Deepfake video of Volodymyr Zelensky surrendering surfaces on social media

**Manipulated images**

A manipulated video of Volodymyr Zelensky has surfaced on social media

DeepFakes of President Zelensky telling his soldiers to surrender released by Russian hackers during the Ukraine war

# BENIGN USES OF DEEPFAKE

- Anonymize voice and faces for privacy

- Cost-effective entertainment production

- Amplify the reach of public messages through language localization

- Bringing back the loved ones or imagining different stages of life



Multilingual social campaign video created using DeepFakes

# RESPONSES FROM TECH COMPANIES

| Company | Policy | Action |
| --- | --- | --- |
| Meta | Removes content that has been manipulated in order to mislead users | Developed a detection tool that reverse engineers a single DF image to its generative model, created Deepfake Detection Challenge (DFDC) in 2020 and open-sourced their dataset |
| Google (YouTube) | Bans manipulated media under scam policies | Released large DeepFake dataset on FaceForensics |
| Twitter | Removes content that deceptively share synthetic or manipulated media that are likely to cause harm | Introduced a "three-pronged test" (3 human answerable questions) to determine if media violates Twitter's policy |
| TikTok | Bans digital forgeries that mislead users by distorting truth and cause harm | Enforces identity check before allowing users to use their face swap filters to prevent unconsented DeepFakes |
| Reddit | Does not allow content that impersonates entities misleadingly | |

# RELATED NEWS

1 Sep 2020
Microsoft (Responsible AI) developed Microsoft Video Authenticator which analyses the percentage chance of a photo or video is artificially manipulated

16 Oct 2020
Qualcomm developed a feature which adds digital signature to each photo taken from Qualcomm Snapdragon smartphones to prove its authenticity when posted online

27 Jan 2022
Coalition for Content Provenance and Authenticity (C2PA) had partnered with tech giants including Microsoft, Intel, and Adobe to combat the rapid spread of DeepFake

29 Apr 2022
AI Singapore (AISG) hosted a DeepFake detection competition Trusted Media Challenge. The winner aims to incorporate his AI model into ByteDance's BytePlus platform to make it available to users

Trusted Media Challenge: https://arxiv.org/pdf/2201.04788v2.pdf

# DEEPFAKE GENERATION TYPES

Identity swap



Expression re-enactment



Attribute manipulation

Entire face synthesis

# DEEPFACELAB

Most popular DeepFake generation tool

# DEEPFAKE WORKFLOW

## Data Processing

Curate a database of source & target faces

## Model Training

Learn the features of both faces

## Merging

Transfer source features to target face

# DATA PROCESSING OVERVIEW


1. Data collection


2. Frame extraction


3. Face extraction


4. Landmarking


5. Filtering


6. Masking

# 3. FACE EXTRACTION

Using Single Shot Scale-invariant Face Detector (S3FD)

Red box = head

Blue box = whole face

Grey area = full face

Green lines = landmarks

# 3. FACE EXTRACTION

S3FD: Good at detecting frontal faces, weaker at side faces (about 3% error rate)

# 4. LANDMARKING

S3FD: Relatively small amount of landmarks (30~40), unable to capture rich emotions

Real

Fake

# 4. LANDMARKING

Google's Firebase ML Kit with more detailed landmarks

# 5. SORTING

Goal: cover as many expressions, lighting, and angle as possible

Challenge: remove duplicate faces from a pool of very similar images



Highly similar faces are removed using Visipics



Tip: keep both faces!

# 5. SORTING

Lack of training data results in pitch black mouth area

# 6. MASKING

Remove unrelated features such as hair, glasses, hand, etc.

Generic pre-trained XSeg model performs well in most cases

# 6. MASKING



Detect texture inconsistency

# MODEL TRAINING

# MODEL TRAINING



source

# MODEL TRAINING



fake source

# MODEL TRAINING



target

# MODEL TRAINING



fake target

# MODEL TRAINING



fake source + target

# MODEL TRAINING

Let the model recognize generic human faces first, then specialize on our source & target's faces

Hyperparameters:

| Realistic aspect | Imaginative aspect |
|---|---|
| Eyes & mouth priority | GAN power |
| Uniform yaw | Face style power |
| Learning rate dropout | Background style power |
| Random warp | |
| Masked training | |

# MERGING

Use our trained model to transform face from target to source

Requires manual fine-tune

# DEEPFAKE DETECTION

Essentially an evaluation on a DeepFake's production quality

Cryptography concept: look for easy-to-evaluate but hard-to-forge features

Examples:
- Frame level smoothness
- Unnatural movement/expression
- Inconsistencies near masking area
- Reflection in the eyes
- Phoneme-viseme mismatches

# STATE-OF-THE-ART DETECTION

Multimodal approach: predict the authenticity of video, audio, lip sync (video + audio)

Hashing approach: check against hashes of known examples of in-the-wild DeepFakes

# THANK YOU

# DFL: FACE EXTRACTION

(a) **Heatmap-based facial landmark algorithm 2DFAN** (for faces with standard pose)
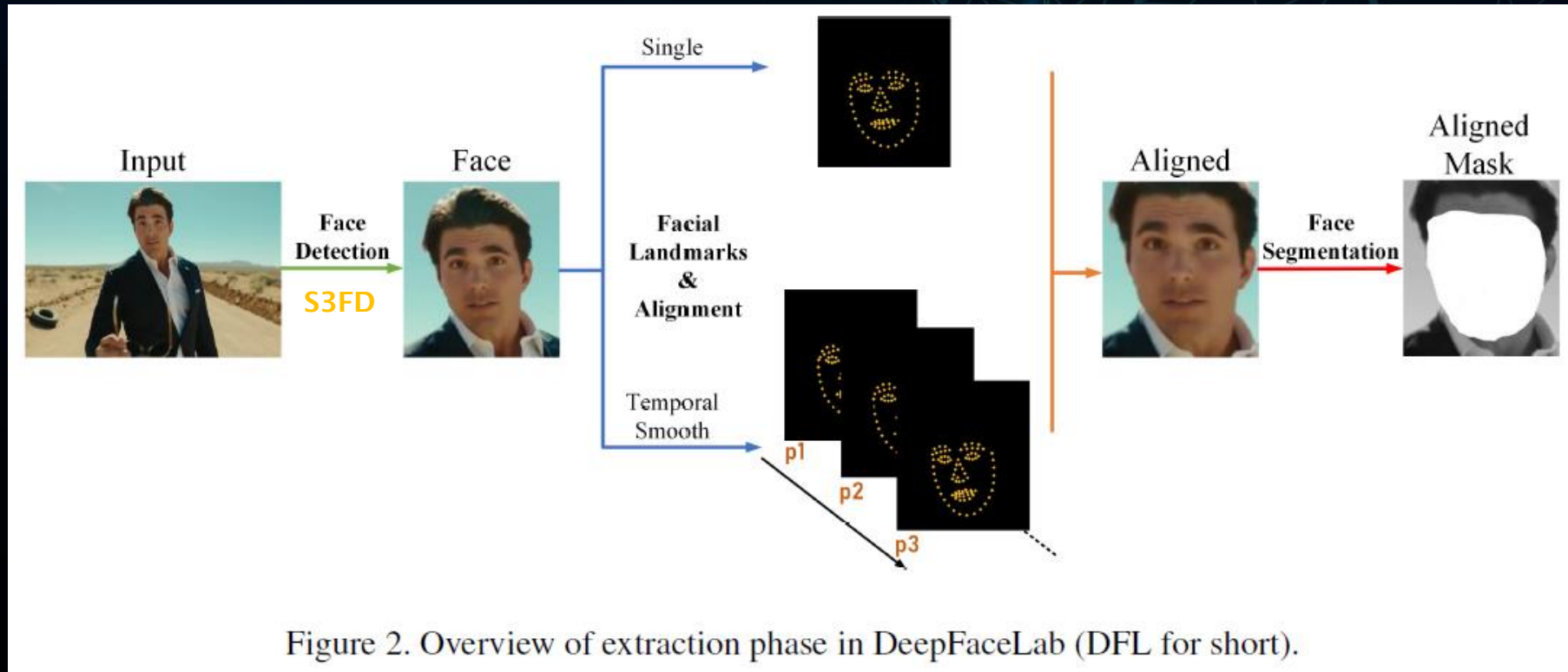(b) **PRNet with 3D face prior information** (for faces with large Euler angle – yaw, pitch, roll)



Figure 2. Overview of extraction phase in DeepFaceLab (DFL for short).

# DFL: MODEL STRUCTURE

Loss:
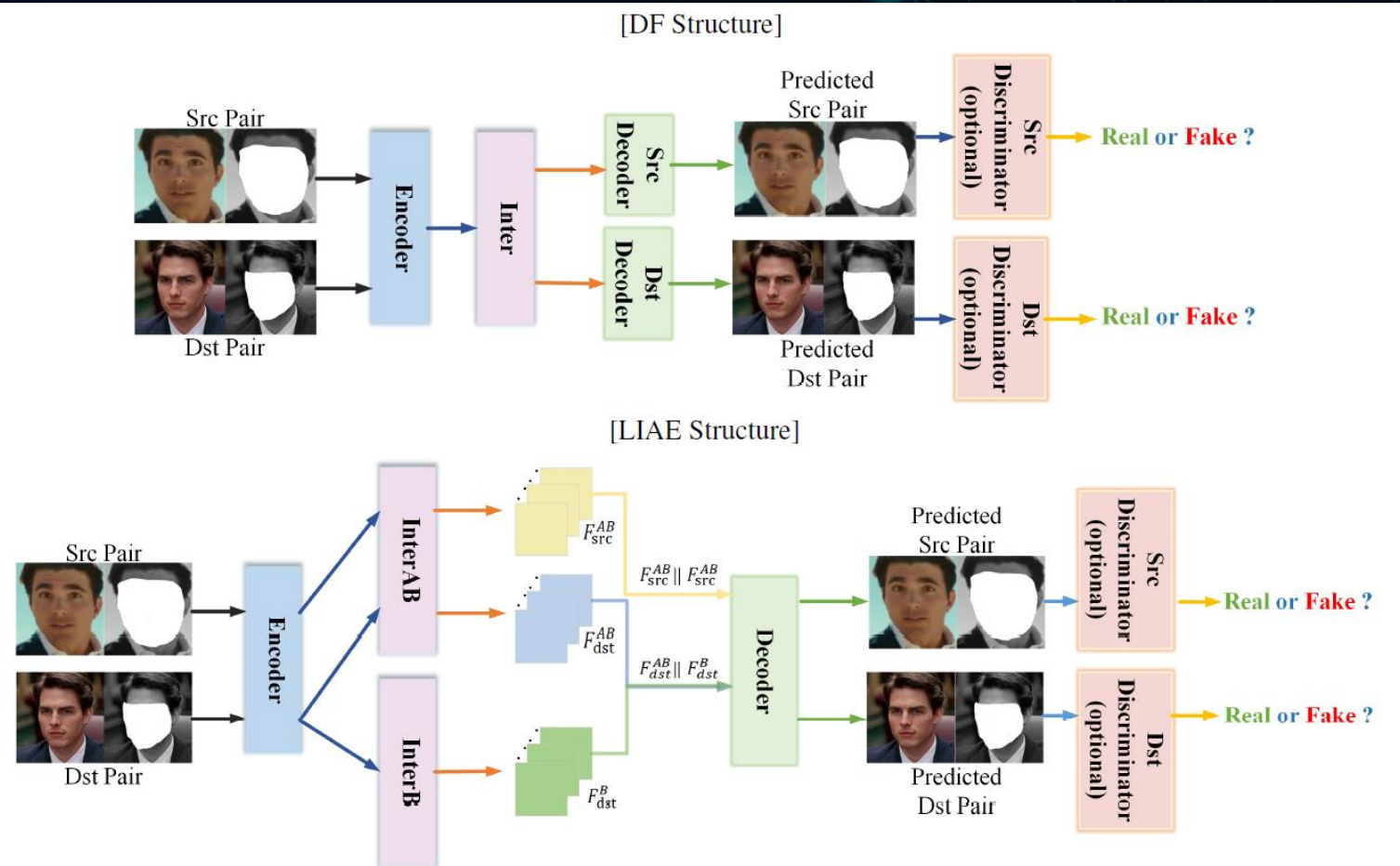(1) DSSIM (structural dissimilarity): faster face generalization
(2) MSE: better clarity



Figure 3. Overview of training phase in DeepFaceLab (DFL). DF structure and LIAE structure are both provided here for illustration, ○||○ represents the concatenation of latent vectors.
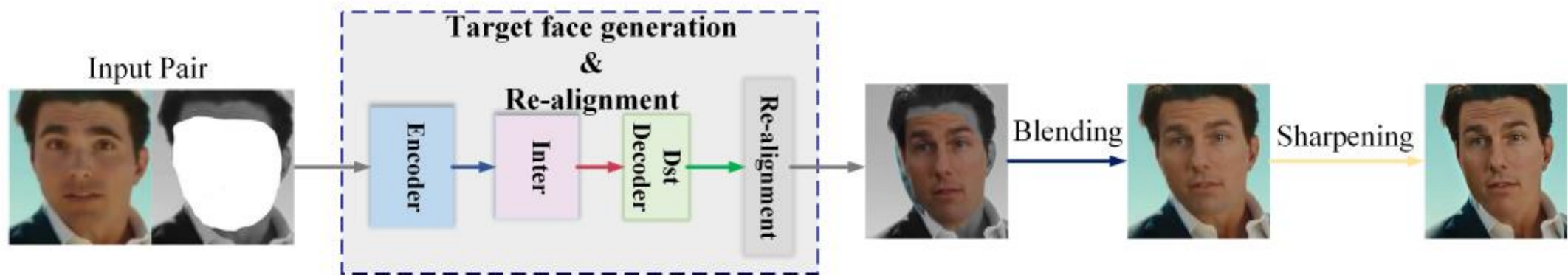
# DFL: MERGING



Figure 4. Overview of conversion phase in DeepFaceLab(DFL).