

# Deepfake Retrieval Systems: Detecting Identity Fraud in Image Databases

---

Jumin Lee and Suhyeon Ha (T4)

2024. 06. 05.

## Introduction

---

# Target Task

- Given an authentic image, our goal is to detect fake images pretending to depict the same person in database.

ID 13, Real



Query

Real & fake images of multiple IDs



Database

ID 13, Fake



Results

# Our Framework

---

- No existing work to retrieve deepfakes of the query image.
- A combination of face retrieval and forgery detection can be utilized.

Stage 1	Stage 2
Face Retrieval	Forgery Detection

or

Stage 1	Stage 2
Forgery Detection	Face Retrieval

- Face retrieval  
: Identify images that match the given identity.
- Forgery detection  
: Determine whether the identified **arbitrary images** have been manipulated.

# Our Framework

---

- Prompt-guided inpainting can modify images while preserving their identities.
- If we use deepfake detection instead of forgery detection, we can not handle this issue.

Two men in a wedding



Two men in jail



Prompt

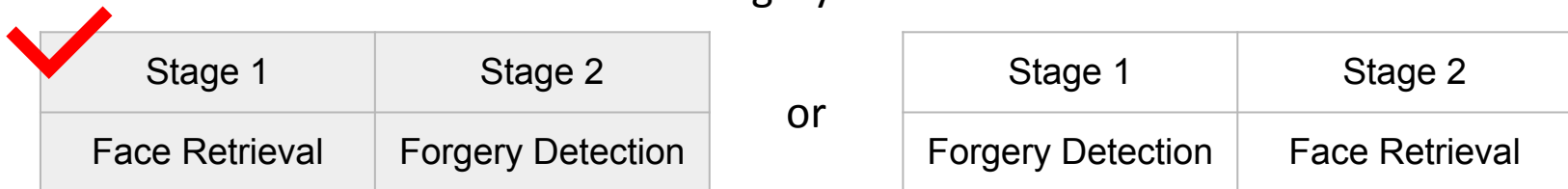
Source Image

Edited Image

# Our Framework

---

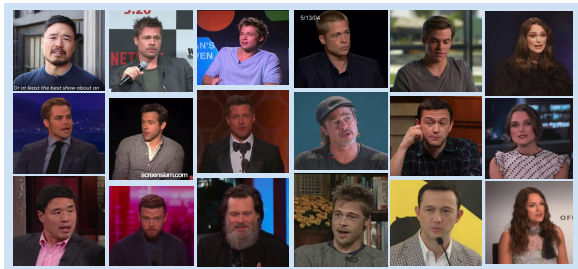
- No existing work to retrieve deepfakes of the query image.
- A combination of face retrieval and forgery detection can be utilized.



- Why is face retrieval ahead of forgery detection?  
: Being unrecognized as someone's identity suggests its quality is doubtful.

# Our Framework

- Given an authentic image, our goal is to detect fake images pretending to depict the same person in database.



Database



ID 13, Real

Stage 1.  
Face Retrieval



ID 13,  
Real&Fake

Stage 2.  
Forgery  
Detection



ID 13, Fake

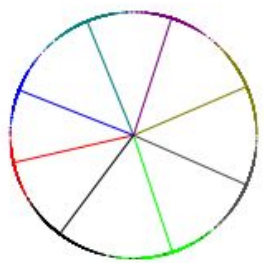
## Related Work

---

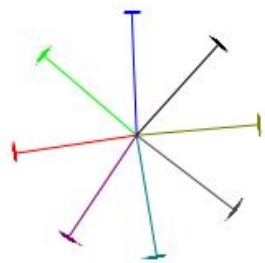


# ArcFace

- ArcFace: Additive Angular Margin Loss for Deep Face Recognition, CVPR 2019



(a) Softmax



(b) ArcFace

< Training >



Detect & Crop



ArcFace



ID  
embedding

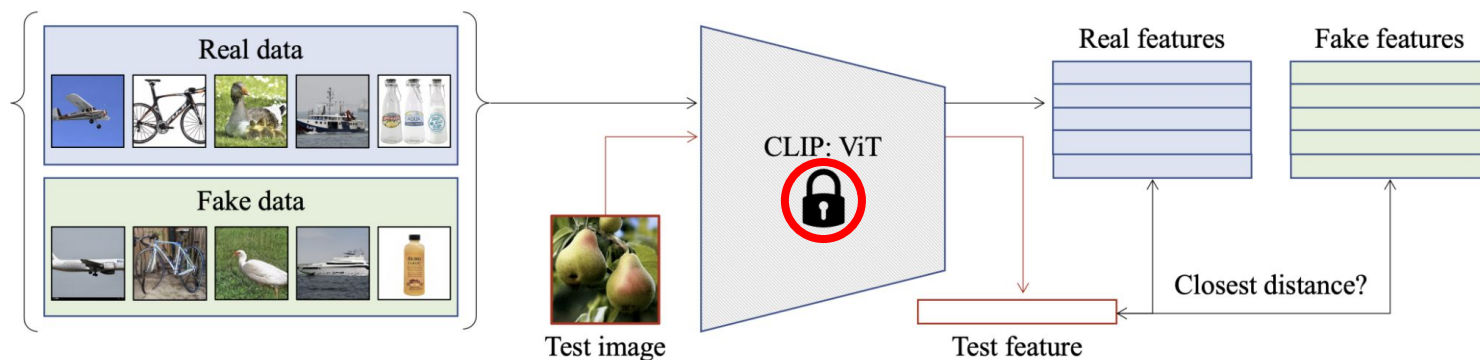


< Inference >

- Train with loss term that depends on the angle between classes to create a larger gap between different classes.
- After training, model can get ID embedding.

# UniDet

- Towards Universal Fake Image Detectors that Generalize Across Generative Models, CVPR 23

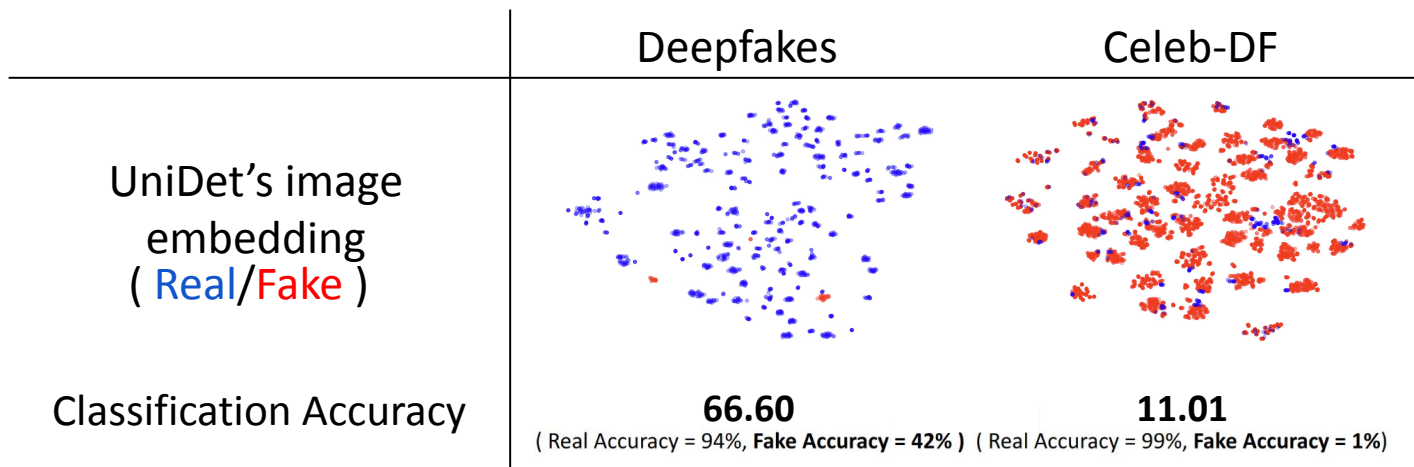


The classification process should happen in a **feature space** which has not been trained to separate images from the two classes.

# In Midterm Project Presentation

- Checked the performance of UniDet for stage 2.
- Classification accuracy result.

Detection method	Generative Adversarial Networks						Denoising Diffusion Models			DALL-E	Deepfake Models	
	Pro-GAN	Cycle-GAN	Big-GAN	Style-GAN	Gau-GAN	Star-GAN	Glide	Guided	LDM		Deepfakes	CelebDF
UniDet w/ LC	100.00	<b>98.50</b>	<b>94.50</b>	82.00	<b>99.50</b>	97.00	79.07	70.03	<b>94.19</b>	81.47	66.60	11.01

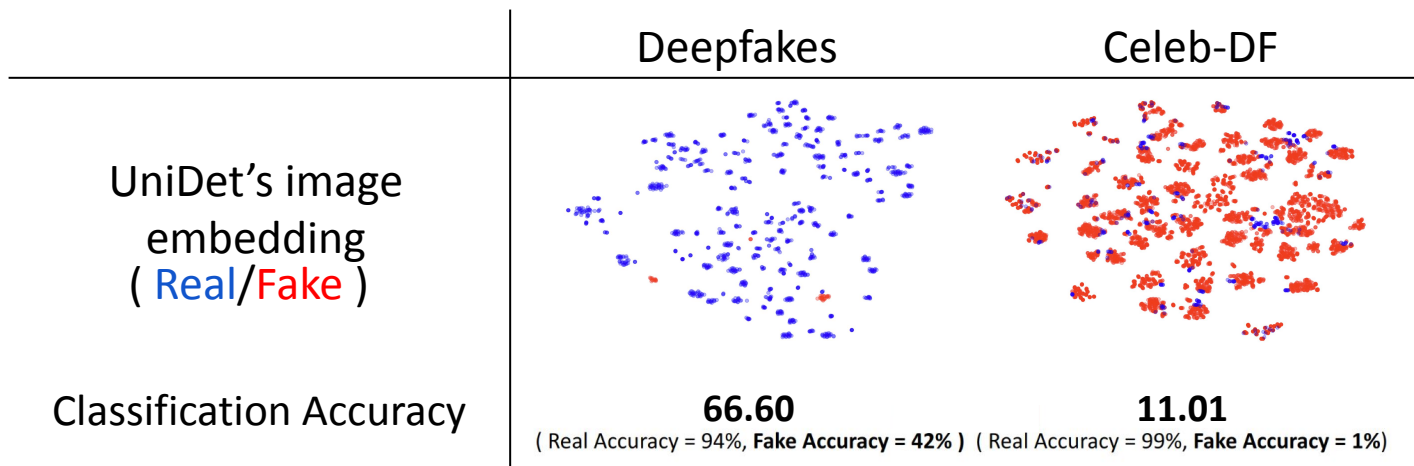


# In Midterm Project Presentation

- There is still room for improvement in forgery detection for both datasets.

→ Our goal is to improve facial forgery detection of UniDet for deepfake retrieval system.

Detection method	Generative Adversarial Networks						Denoising Diffusion Models			DALL-E	Deepfake Models	
	Pro-GAN	Cycle-GAN	Big-GAN	Style-GAN	Gau-GAN	Star-GAN	Glide	Guided	LDM		Deepfakes	CelebDF
UniDet w/ LC	100.00	<b>98.50</b>	<b>94.50</b>	82.00	<b>99.50</b>	97.00	79.07	70.03	<b>94.19</b>	81.47	66.60	11.01



# Our Approach

---

# Our Approach

---

- UniDet uses only CLIP's visual features for forgery detection.
- We try to combine CLIP's visual and **text features** to improve UniDet's performance.
  - It's used for many other tasks such as classification and generation.  
ex ) Classification : CoOp[1], CoCoOp[2]  
ex ) Generation : Arc2Face[3]

[1] Learning to Prompt for Vision-Language Models, IJCV 2022

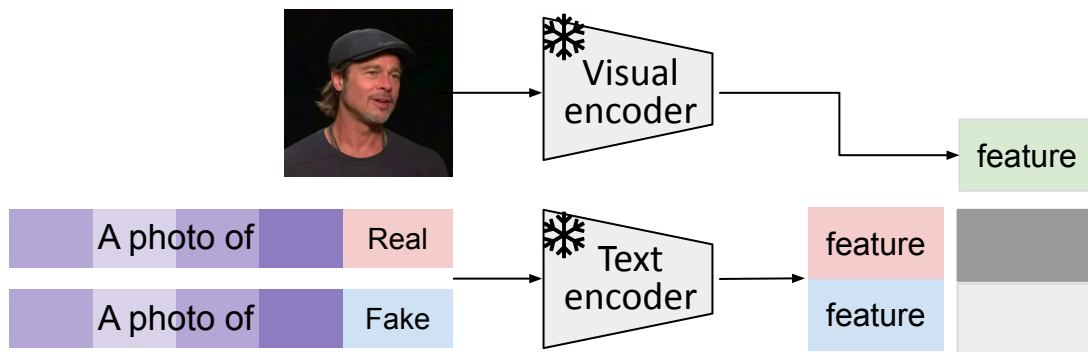
[2] Conditional Prompt Learning for Vision-Language Models, CVPR 2022

[3] Arc2Face: A Foundation Model of Human Faces, arxiv 2024

# Our Approach

---

- UniDet uses only CLIP's visual features for forgery detection.
- We try to combine CLIP's visual and **text features** to improve UniDet's performance.

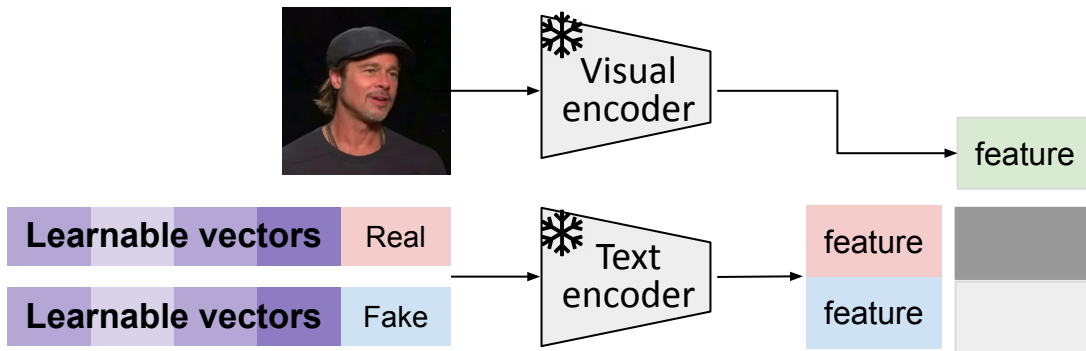


- However, we need prompt engineering which is inefficient.
  - A photo of [CLASS]. / A photo of a [CLASS]. / A [CLASS]. / ...

# Our Approach

---

- UniDet uses only CLIP's visual features for forgery detection.
- We try to combine CLIP's visual and **text features** to improve UniDet's performance.

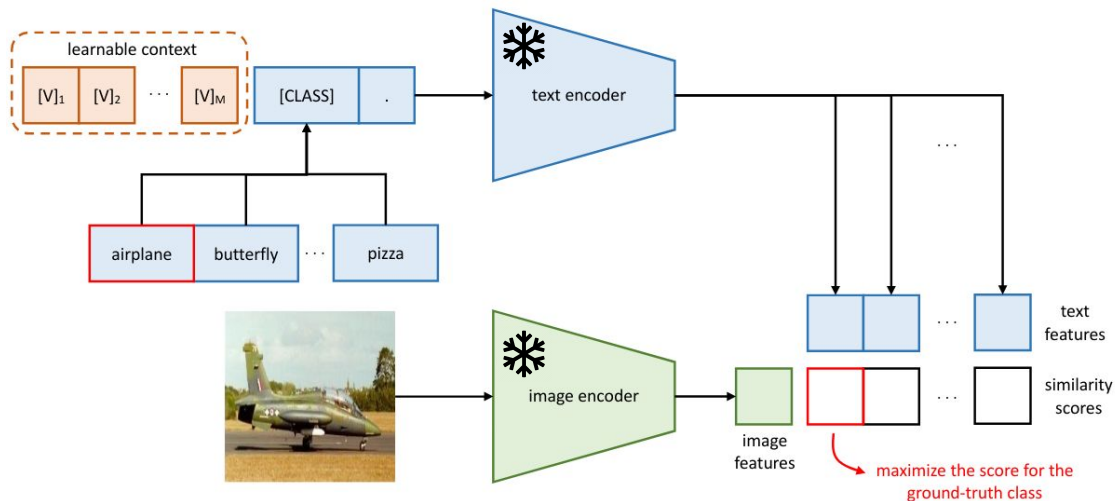


- So we apply Context Optimization (CoOp).



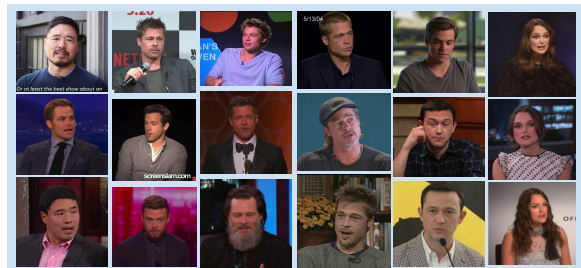
# Our Approach

- What is Context Optimization (CoOp) ?
  - : Model prompt's context words with learnable vectors while the entire pre-trained parameters are kept fixed.
  - Using cross entropy loss.



# Our Framework

- Inference



Database



ID 13, Real

Face Retrieval

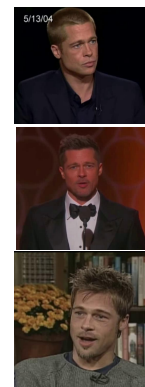
Stage 1.  
ArcFace



ID 13,  
Real&Fake

Forgery  
Detection

Stage 2.  
Ours



ID 13, Fake

## Results

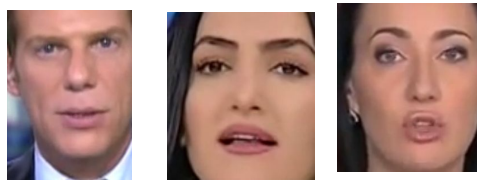
---

# Experiment details for Stage 2.

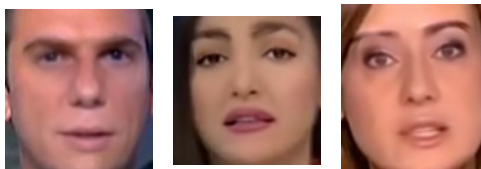
---

- Training Dataset : ProGAN dataset
- Evaluation Dataset for Deepfake Models
  - Deepfakes : 5,405 frames ( = 2,707 real + 2,698 fake)
  - CelebDF : 50,205 frames ( = 4,820 real + 45,385 fake)

Real



Fake



Deepfakes

CelebDF

# Quantitative results

- Stage 2. Classification Accuracy.

Detection method	Generative Adversarial Networks						Denoising Diffusion Models			DALL-E	Deepfake Models		Avg.
	Pro-GAN	Cycle-GAN	Big-GAN	Style-GAN	Gau-GAN	Star-GAN	Glide	Guided	LDM		Deepfakes	CelebDF	
UniDet w/ LC	100.00	<b>98.50</b>	<b>94.50</b>	82.00	<b>99.50</b>	97.00	79.07	70.03	<b>94.19</b>	81.47	66.60	11.01	81.15
Ours	100.00	95.60	93.80	<b>95.25</b>	93.43	<b>99.15</b>	<b>92.88</b>	<b>84.3</b>	88.16	<b>91.50</b>	<b>79.63</b>	<b>11.70</b>	<b>93.22</b>

- Demonstrates **high performance gains**, especially on Deepfakes dataset.
- However, we also have same problem with CelebDF dataset.
  - Even if we crop out just the faces like in the deepfake dataset, we achieve **14.25%** accuracy performance.

# Quantitative results

---

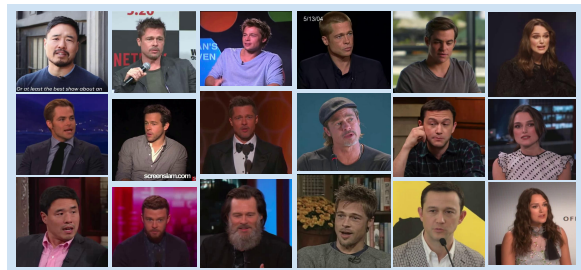
- Stage 2. Classification Accuracy.

Detection method	Deepfake Models	
	Deepfakes	CelebDF
UniDet w/ LC	66.60	11.01
Ours	79.63	11.70
Ours finetuned w/ CelebDF dataset	<b>87.29</b>	<b>50.34</b>

- Fine-tuning with CelebDF datasets can increase the Deepfakes' performance to 87%.
- Still show low performance on CelebDF dataset.

# Our overall framework

- Inference



Database



ID 13, Real

Face Retrieval

Stage 1.  
ArcFace



ID 13,  
Real&Fake

Forgery  
Detection

Stage 2.  
Ours

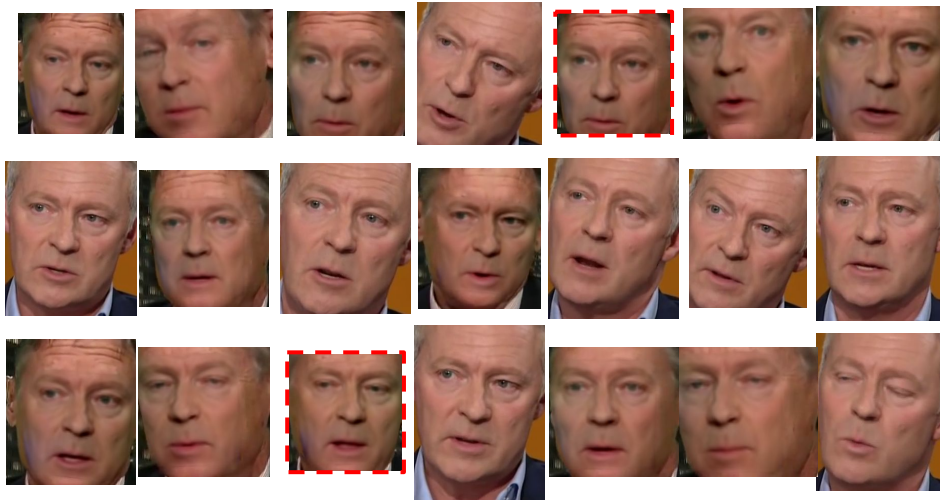


ID 13, Fake

# Qualitative results

Retrieved images from each stage are shown. Red-dotted box denotes missed deepfake.

ID 273, Real



Recall: 100% (22/22)

Query

Stage 1



Recall: 85.7% (12/14)

Stage 2



# Qualitative results

Retrieved images from each stage are shown. Red-dotted box denotes missed deepfake.

ID 731, Real



Recall: 100% (19/19)

Stage 1

Query



Recall: 80% (4/5)

Stage 2

# Qualitative results

- We also check our model can detect modified images while preserving their identities.



Query



Recall: 80% (8/10)

# Conclusion

---

- We first construct the **deepfake retrieval framework**.
  - Not just when identities change, but also when backgrounds change.
- Significant performance improvement compared to UniDet by **using text features**, especially on Deepfakes dataset.
- Limitation
  - The protocol of universal deepfake detection is based on ProGAN, but there is a lack of face images in this dataset, so we need a universal deepfake detection method that can overcome this problem.

Detection method	Generative Adversarial Networks						Denoising Diffusion Models			DALL-E	Deepfake Models		Avg.
	Pro-GAN	Cycle-GAN	Big-GAN	Style-GAN	Gau-GAN	Star-GAN	Glide	Guided	LDM		Deepfakes	CelebDF	
UniDet w/ LC	100.00	<b>98.50</b>	<b>94.50</b>	82.00	<b>99.50</b>	97.00	79.07	70.03	<b>94.19</b>	81.47	66.60	11.01	81.15
Ours	100.00	95.60	93.80	<b>95.25</b>	93.43	<b>99.15</b>	<b>92.88</b>	<b>84.3</b>	88.16	<b>91.50</b>	<b>79.63</b>	<b>11.70</b>	<b>93.22</b>

# Roles

---

- Jumin
  - Implement CoOp
  - Generate PhotoGuard samples
  
- Suhyeon
  - Implement ArcFace
  - Implement inference framework

# Q&A

---

# CS588 Final Project Presentation

---

Thank you.

---