

What's wrong with this video?

Comparing Explainers for Deepfake Detection

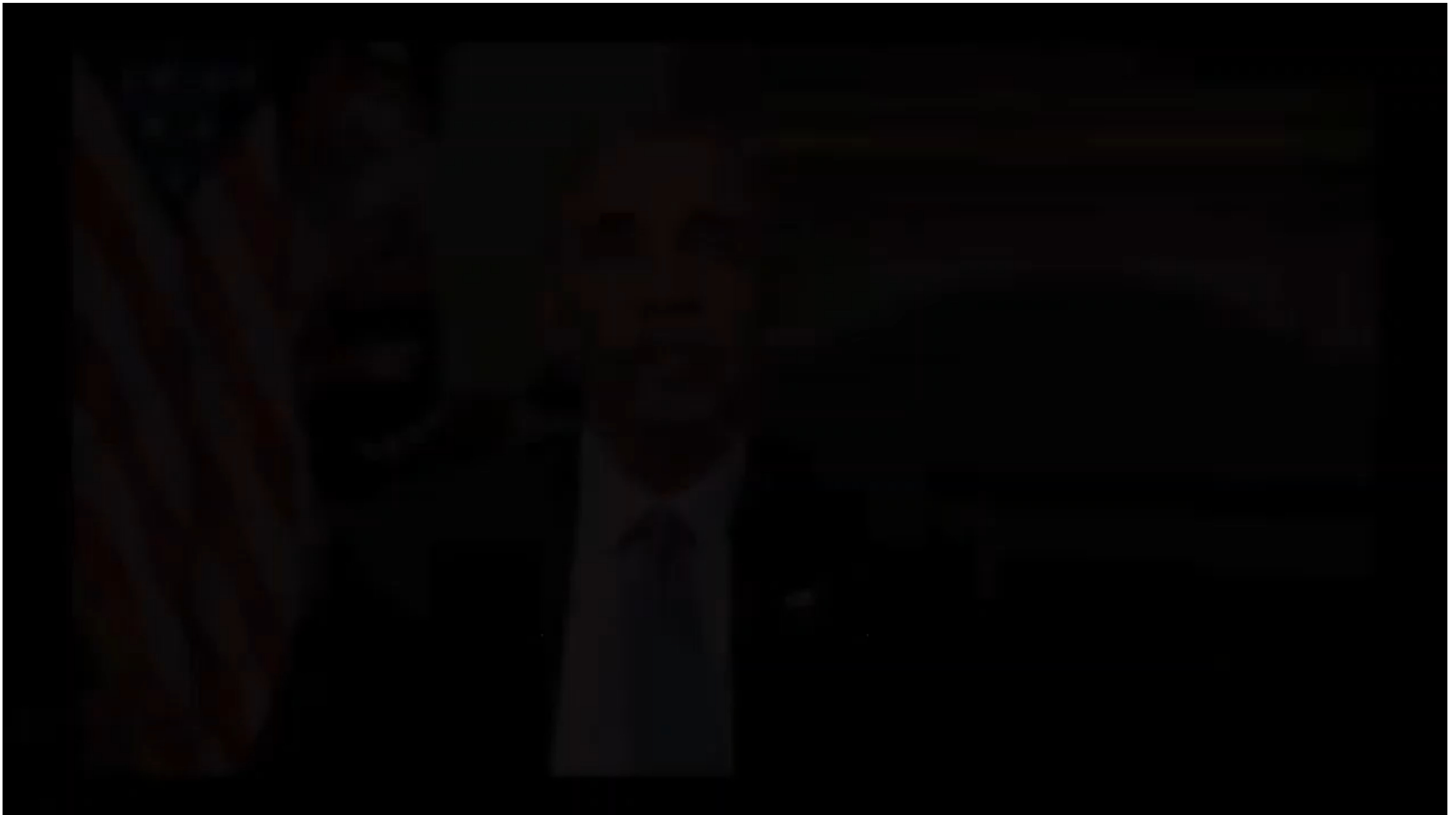
Samuele Pino
samuele.pino@mail.polimi.it
CSE Track



POLITECNICO
MILANO 1863



HP-SR
in Information Technology



<https://www.youtube.com/watch?v=cQ54GDm1eL0>

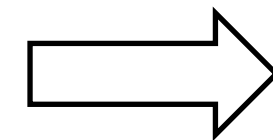
Goals

Automatic classifiers can already detect if a video is real or fake.

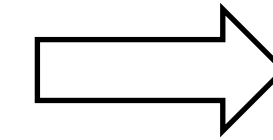
But can we understand the reason why a video is detected as fake?



Video



Manipulation detector



REAL

FAKE

WHY?

In this work we develop, extend and compare explanation techniques for deepfake detection.

Overview

Overview

- Introduction



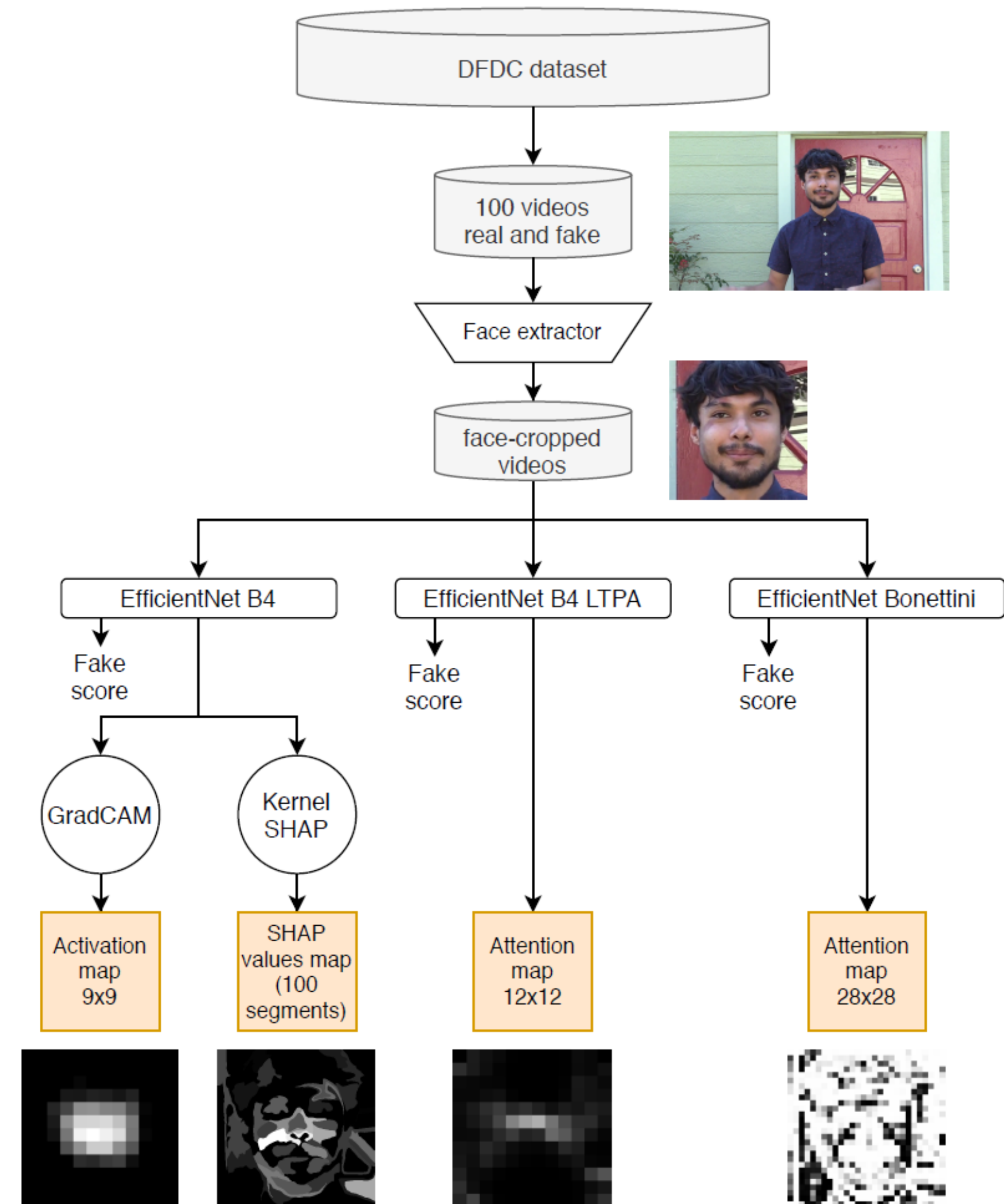
Original showing Alison Brie



Deepfake showing Jim Carrey instead of Brie

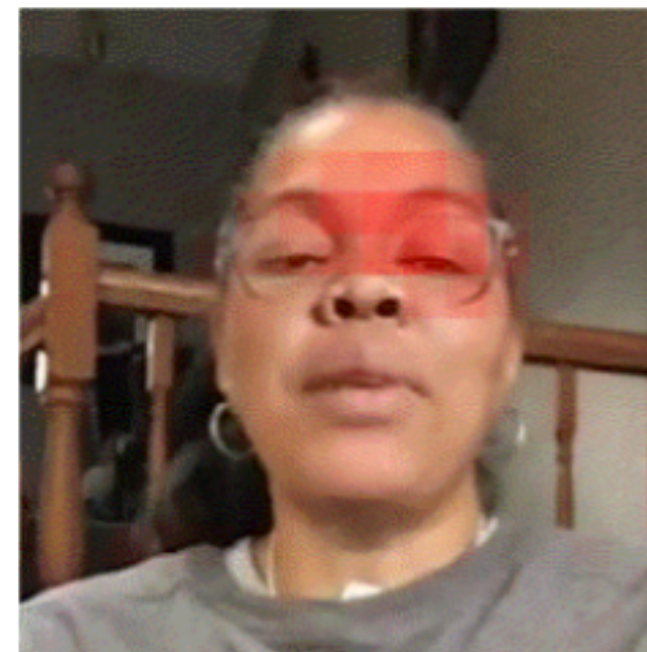
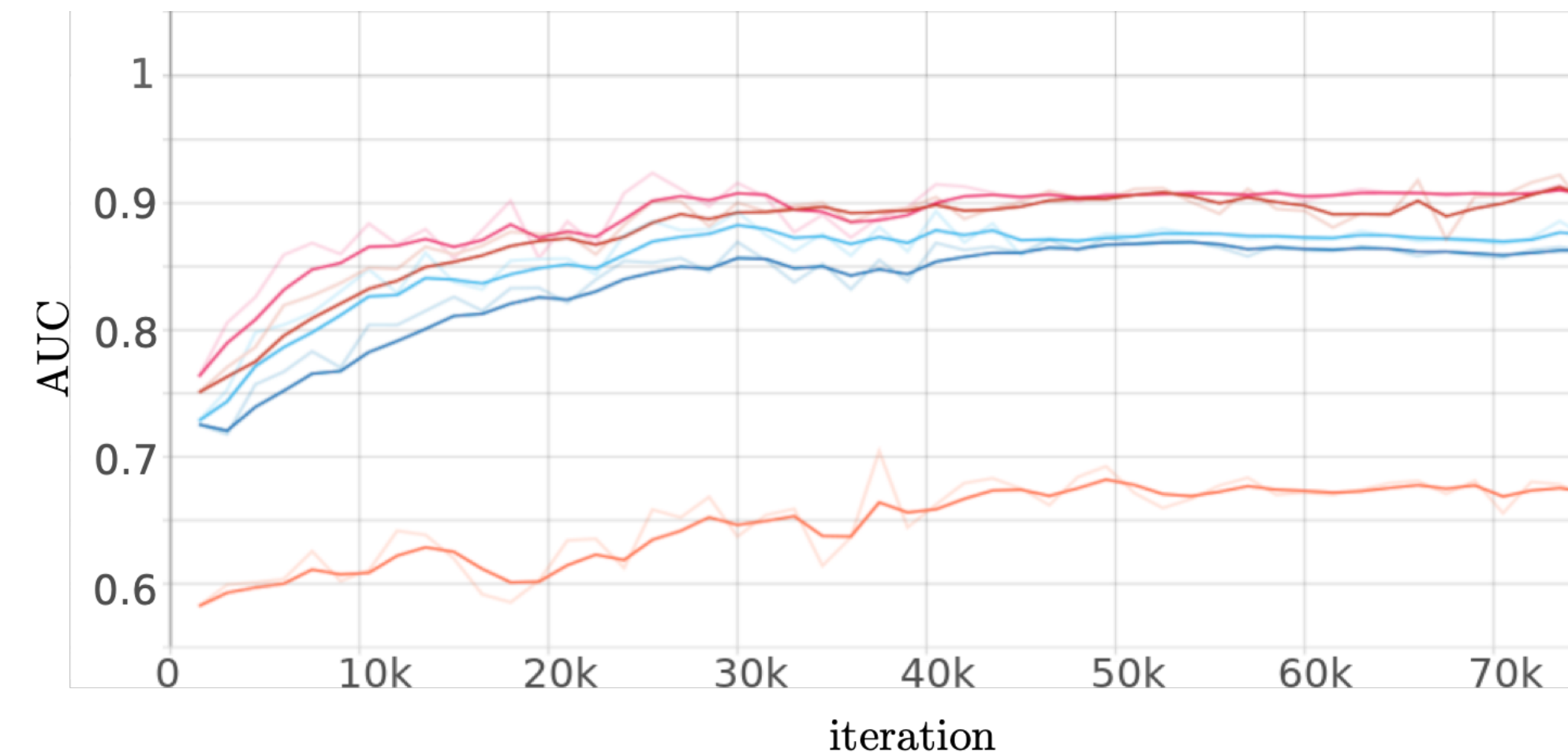
Overview

- Introduction
- Approach

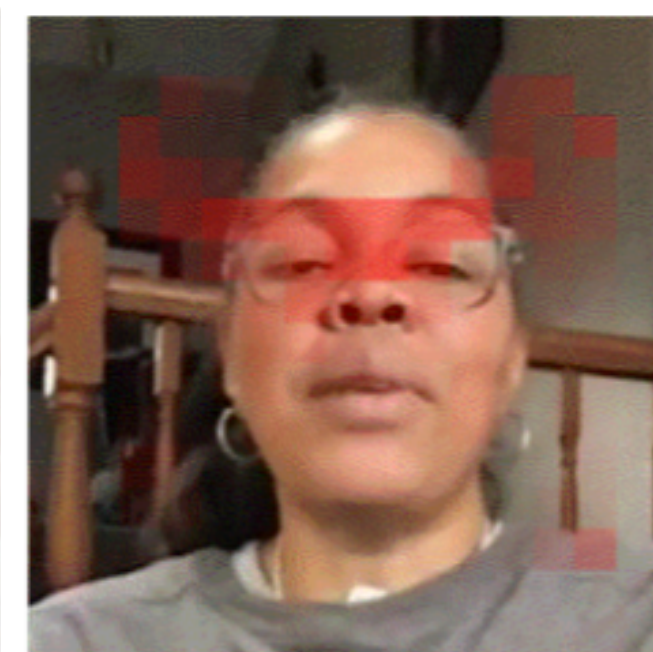


Overview

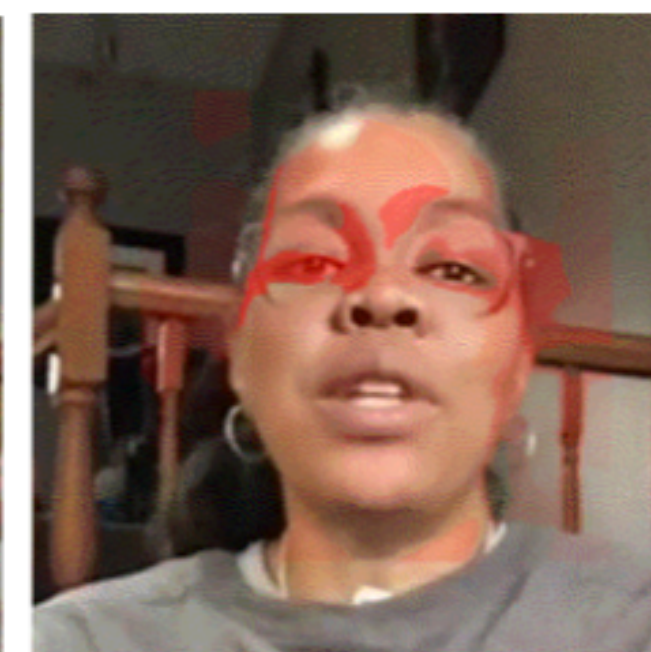
- Introduction
- Approach
- Experiments



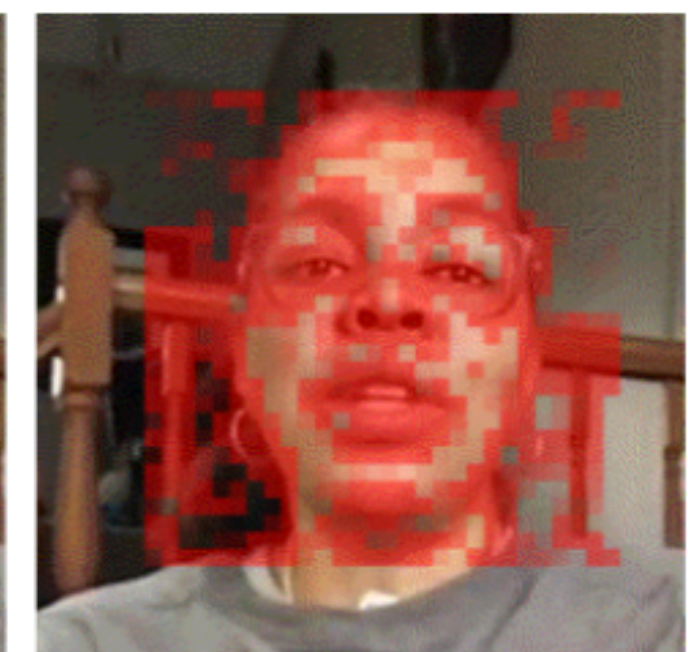
(a) GradCAM



(b) LTPA lv. 2



(c) SHAP



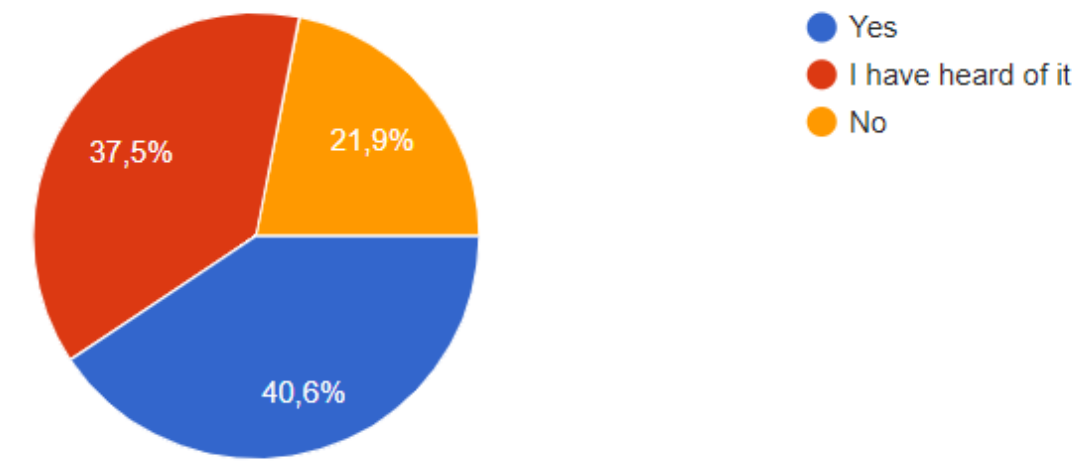
(d) Bonettini

Overview

- Introduction
- Approach
- Experiments
- Results

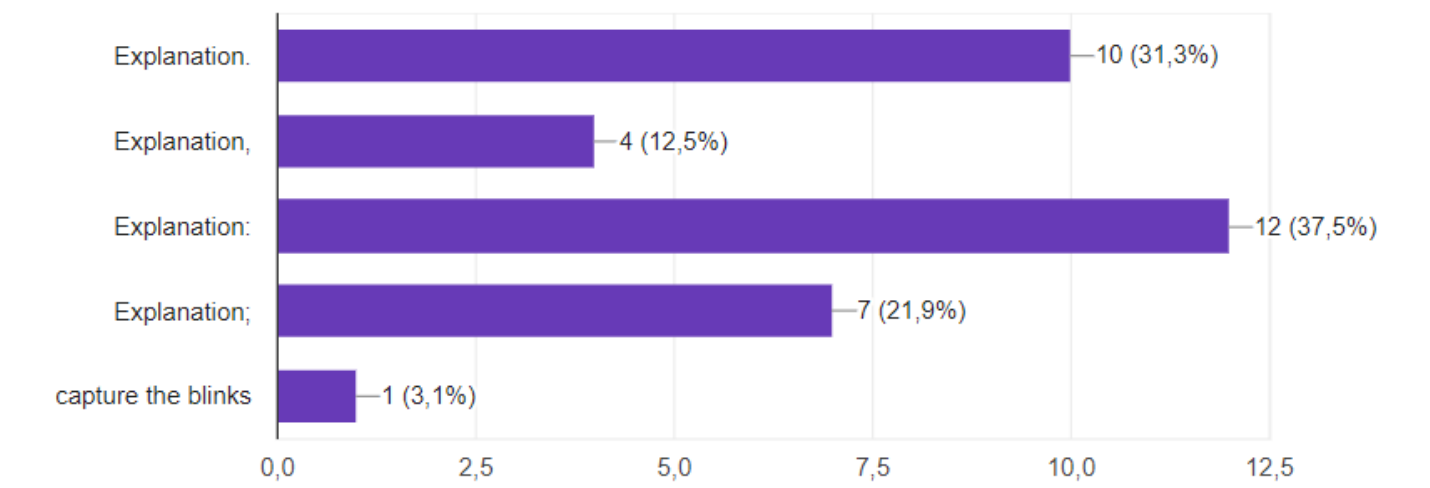
Are you familiar with the concept of deepfake?

32 risposte



Q2: A bot thinks that this face has been edited (indeed it is). In your opinion, which ones of the 4 animations best explain why the robot believes this?

32 risposte



	I [0, 1]	V [0, 1]	τ [-1, 1]	ρ [-1, 1]	μ [0, 1]
Bonettini	0.4821	0.0951	0.7390	0.1262	0.5286
GradCAM	0.0689	0.0135	0.8756	0.7489	0.8666
LTPA	0.0616	0.0108	0.7991	0.3333	0.6386
SHAP	0.0563	0.0302	0.4496	0.2326	0.7348

Introduction

Deepfakes

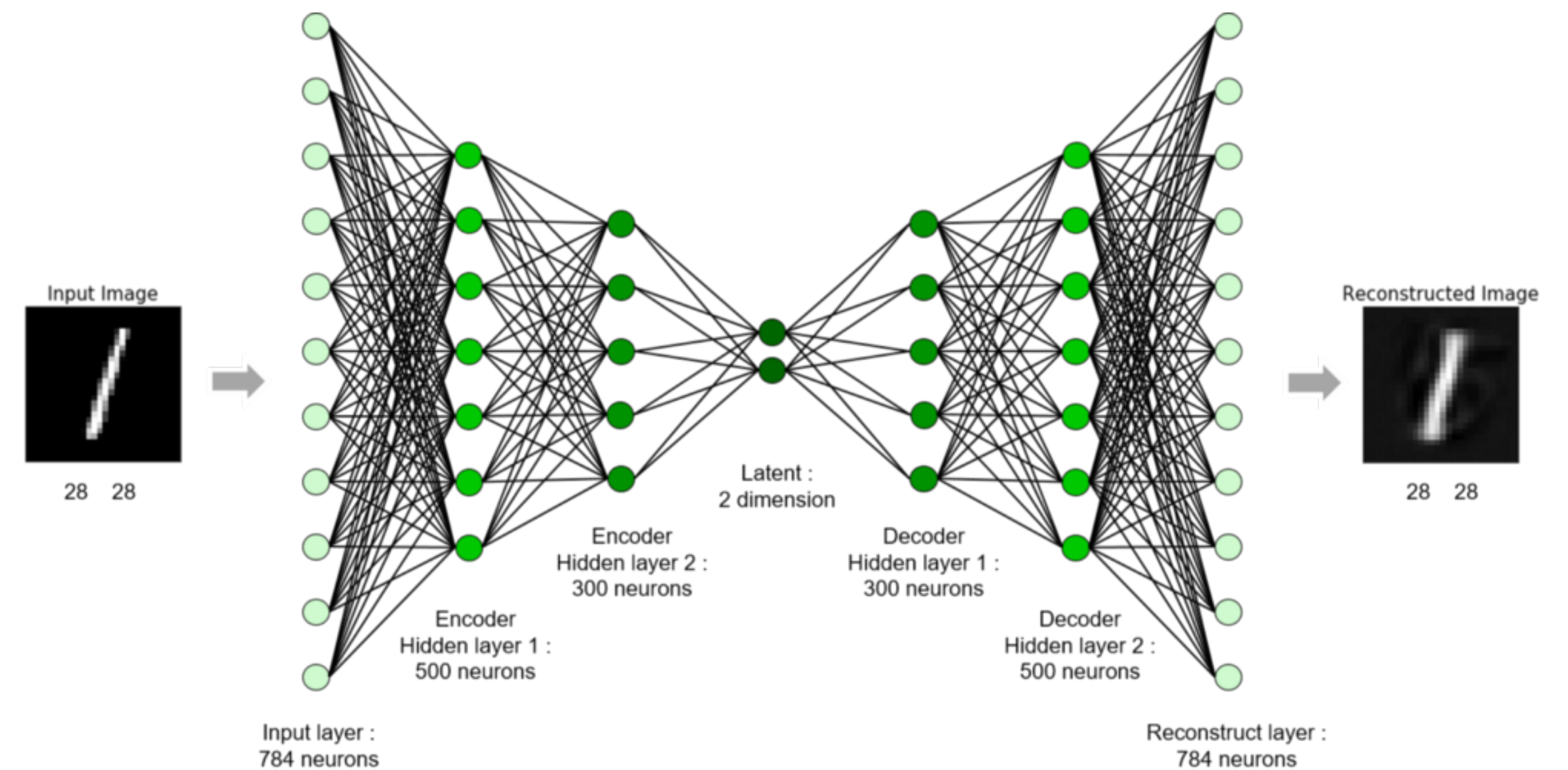
- Replacing faces in videos



[Deepfake Detection Challenge, 2019]

Deepfakes

- Replacing faces in videos
- Deep learning technique



Deepfakes

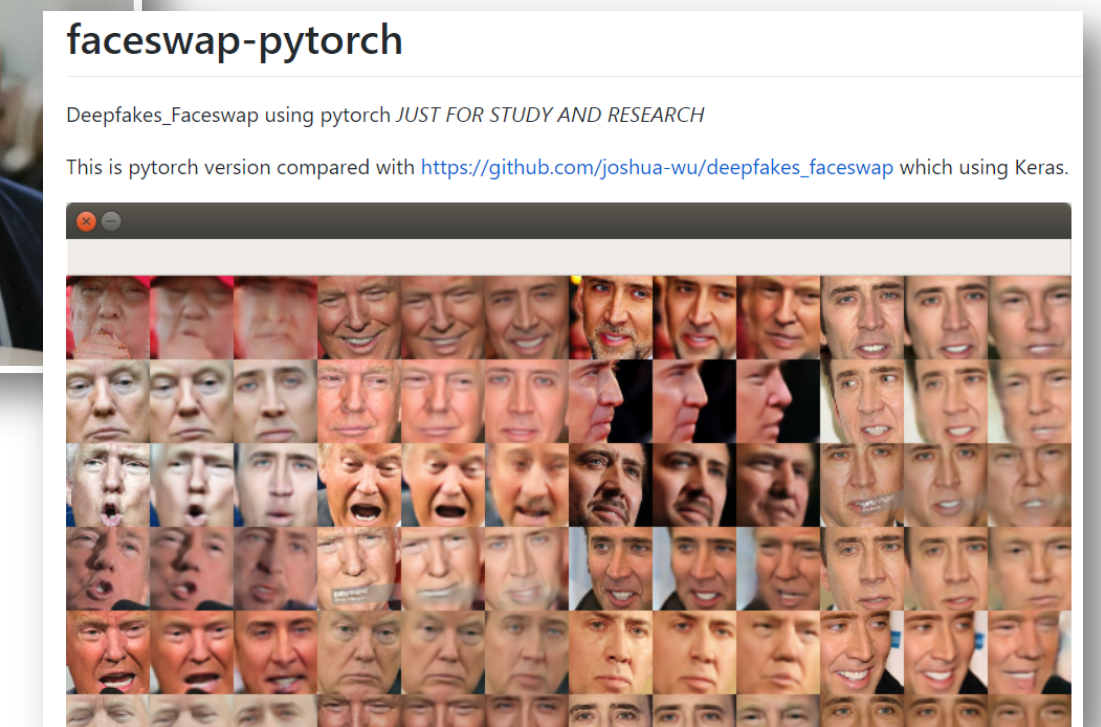
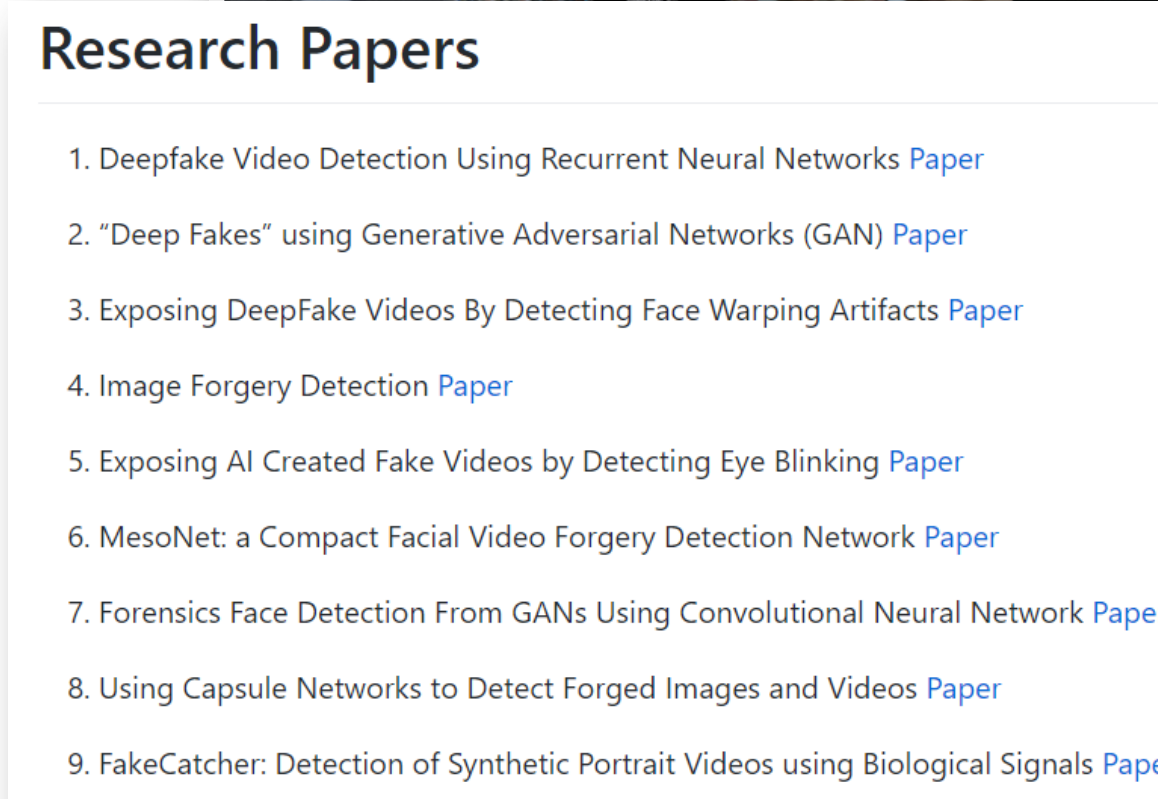
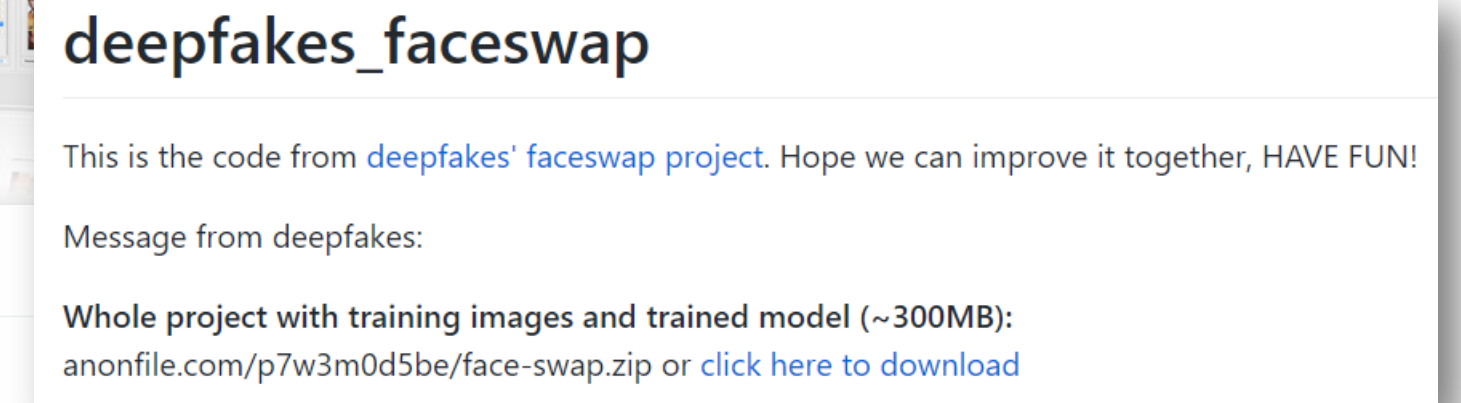
- Replacing faces in videos
- Deep learning technique
- Initially to generate adult contents



r/deepfakes has been banned from Reddit

Deepfakes

- Replacing faces in videos
- Deep learning technique
- Initially to generate adult contents
- No official implementation



Deepfakes

Why is it important to detect them?

Deepfakes

Why is it important to detect them?

- disinformation



Deepfakes

Why is it important to detect them?

- disinformation
- online abuse



Deepfakes

Why is it important to detect them?

- disinformation
- online abuse
- financial fraud



Deepfakes

Why is it important to detect them?

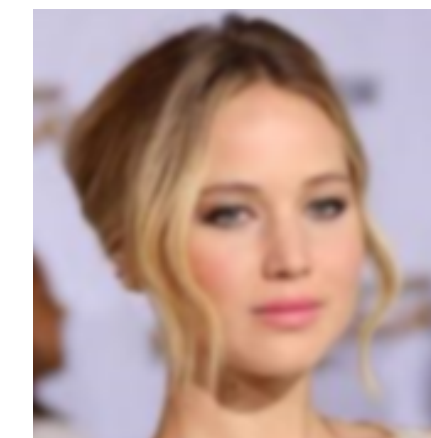
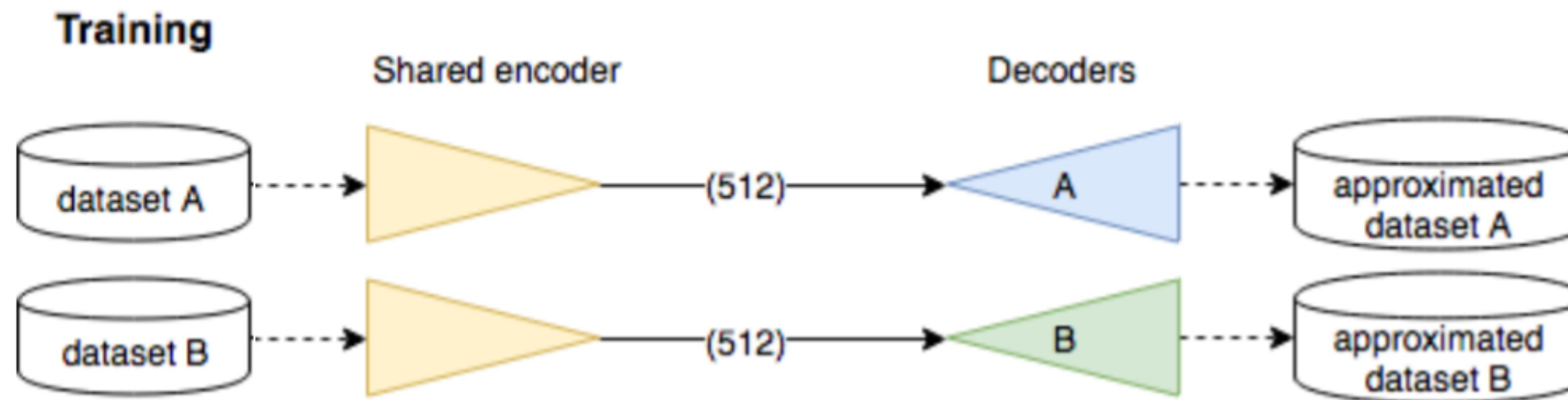
- disinformation
- online abuse
- financial fraud
- law enforcement



Deepfakes

Let's build our deepfake!

Deepfakes

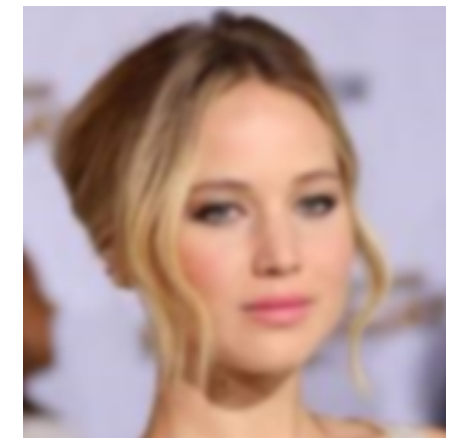


Jennifer Lawrence



Steve Buscemi

Deepfakes



Jennifer Lawrence

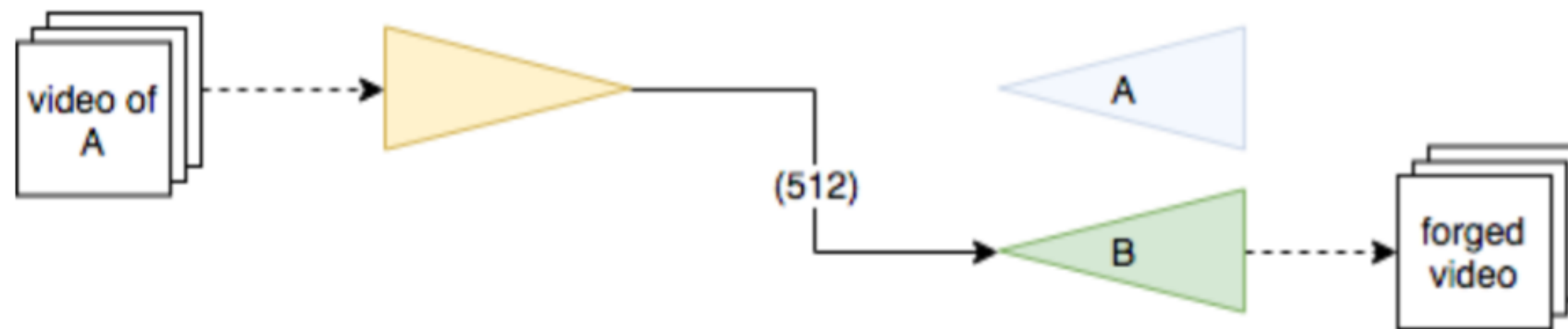


Steve Buscemi

Training



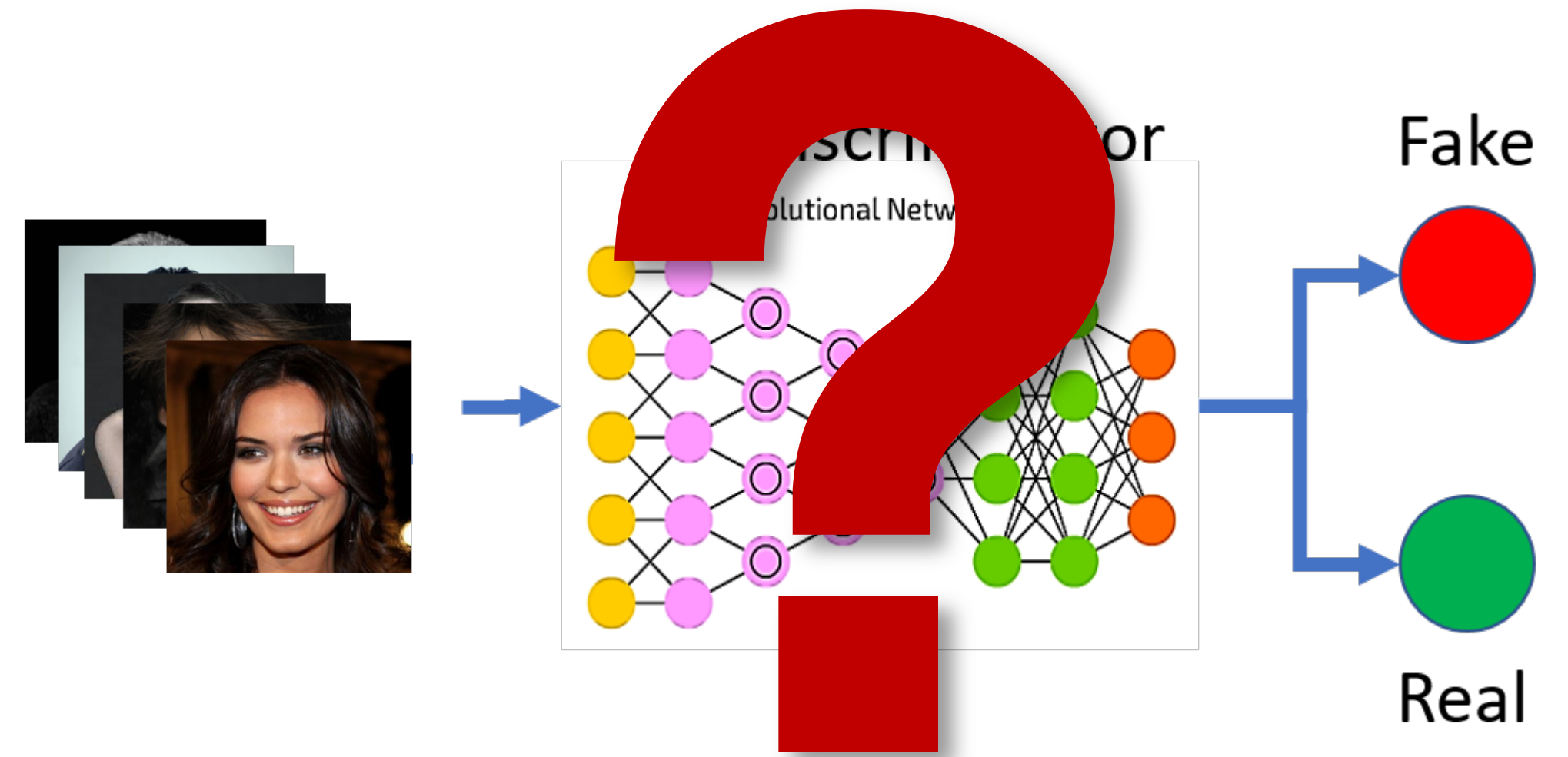
Usage



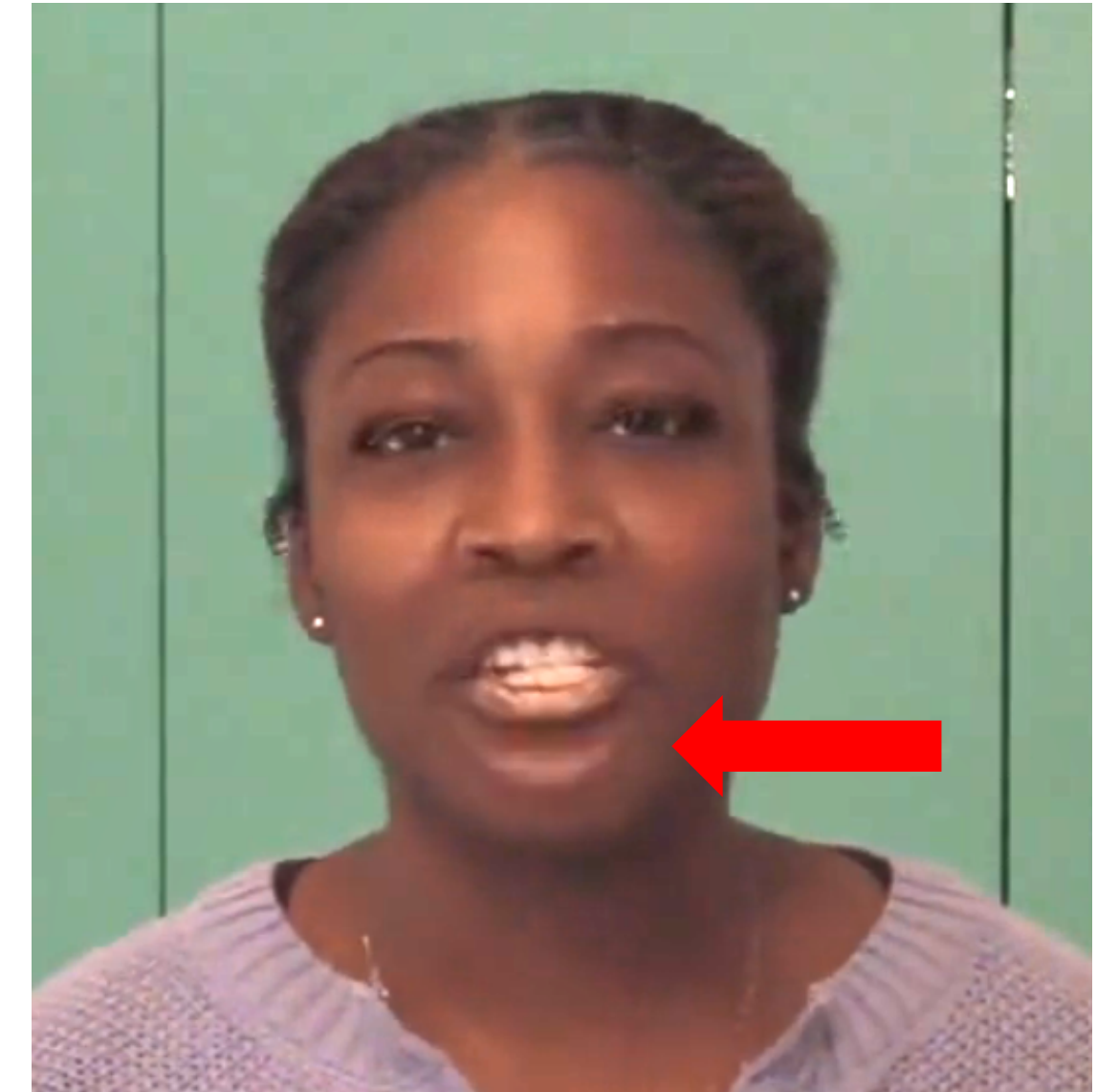
Jennifer Buscemi (?)

Explainability problem

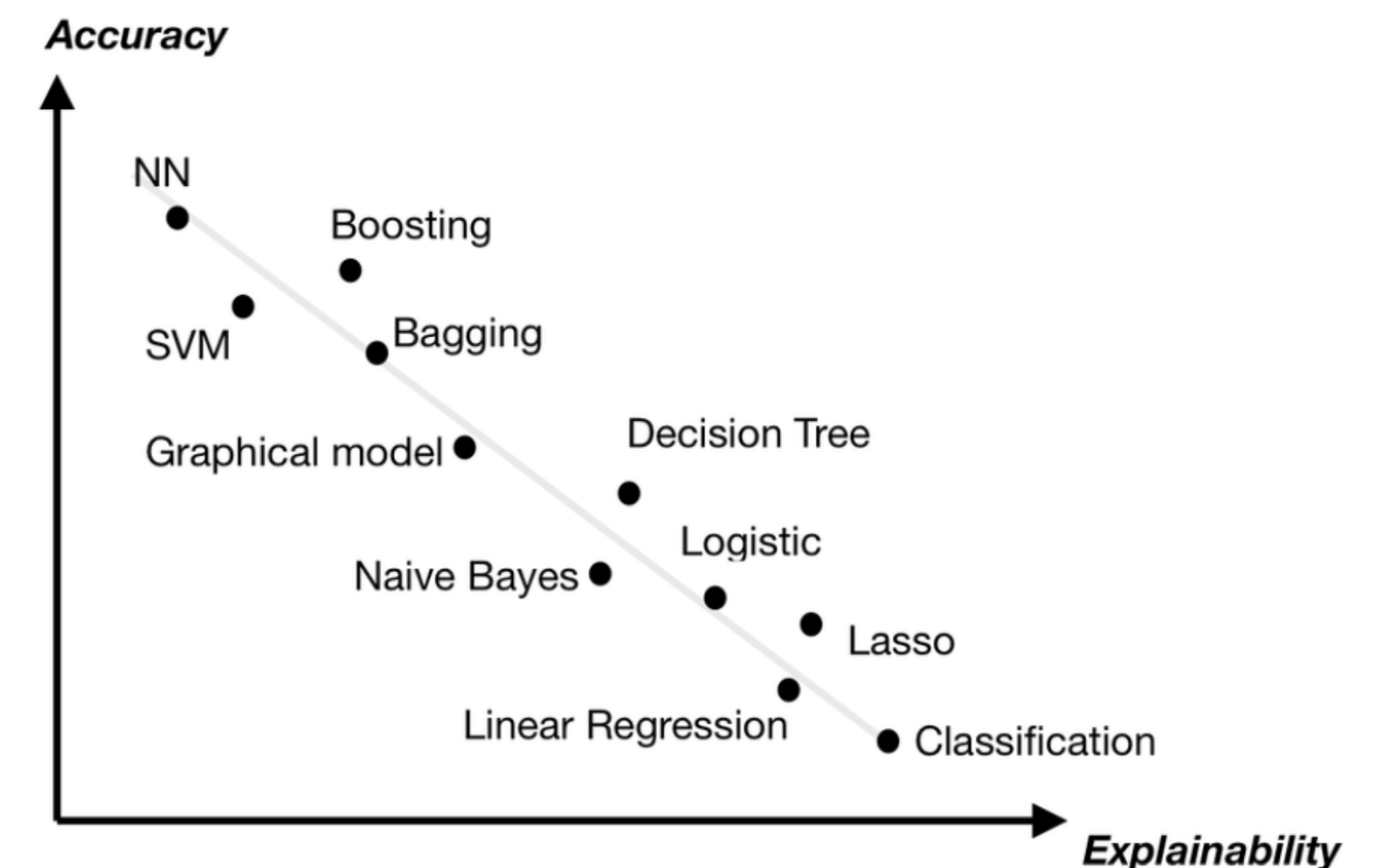
- Also detectors use deep learning
- Should we trust them?



Explainability problem



- «correct prediction for the correct reason»
- Complexity-interpretability trade off



Explainability problem



- Why do we need it:
 - law enforcement
 - journalists
 - dispute resolution in social media



Your video was deleted because it has been detected as fake.

Why would it be fake?!

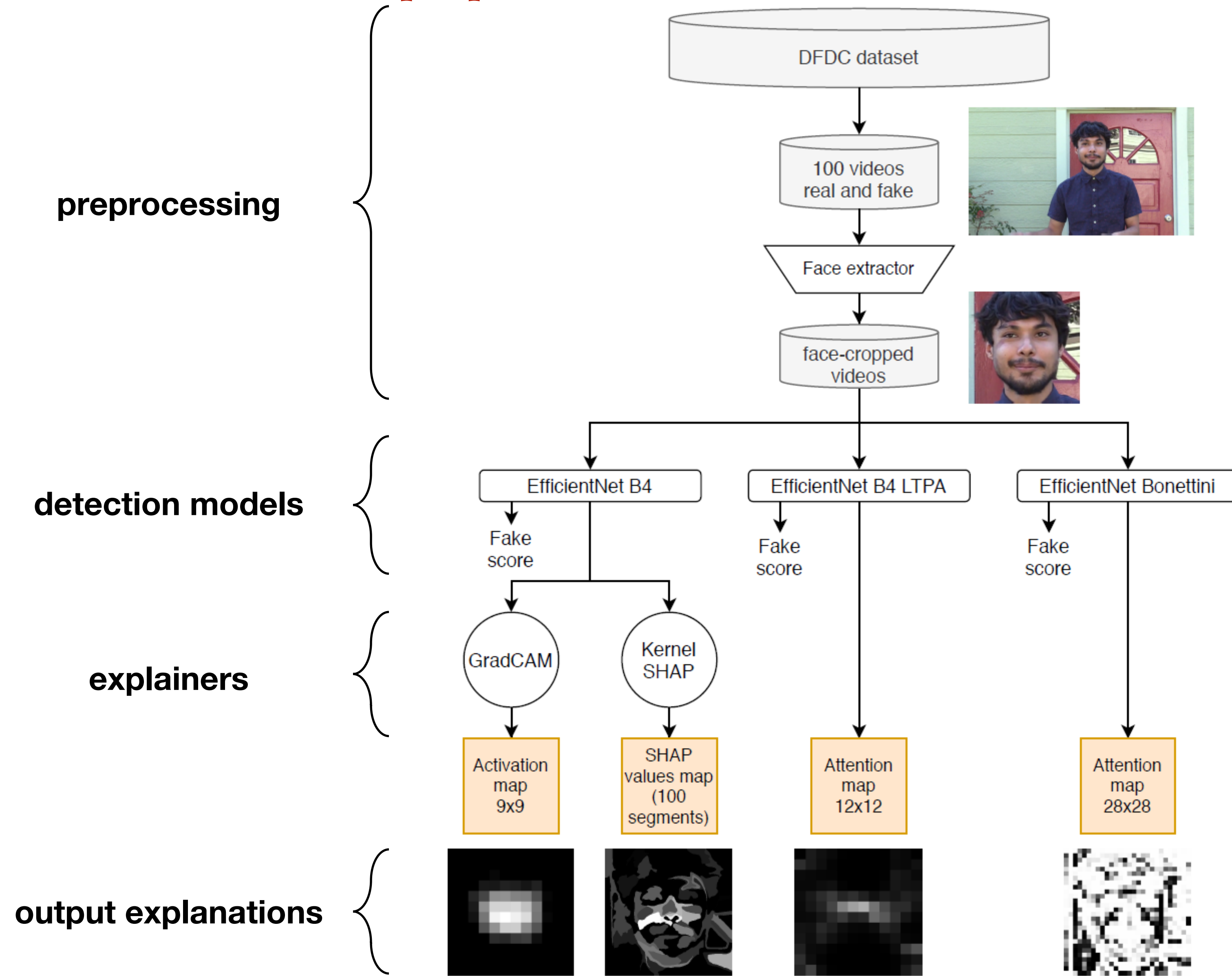


The girl in the video has 2 lower lips.

Ah...

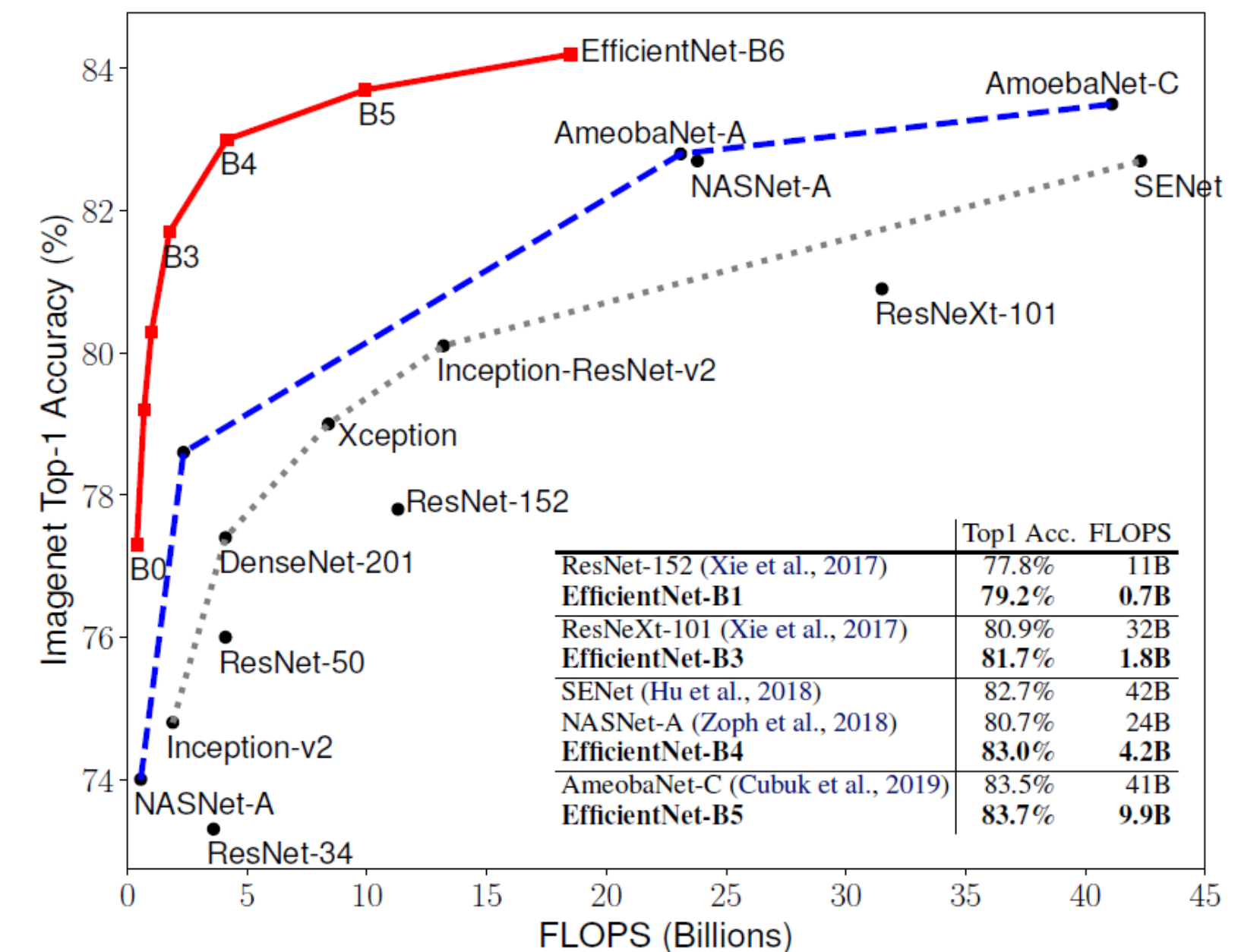
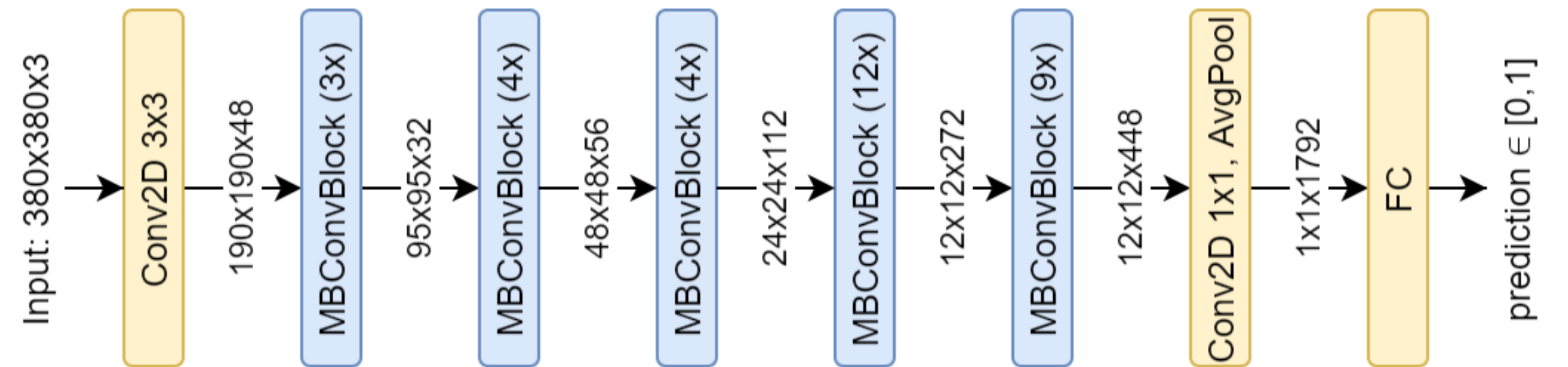
Approach

Approach: overview



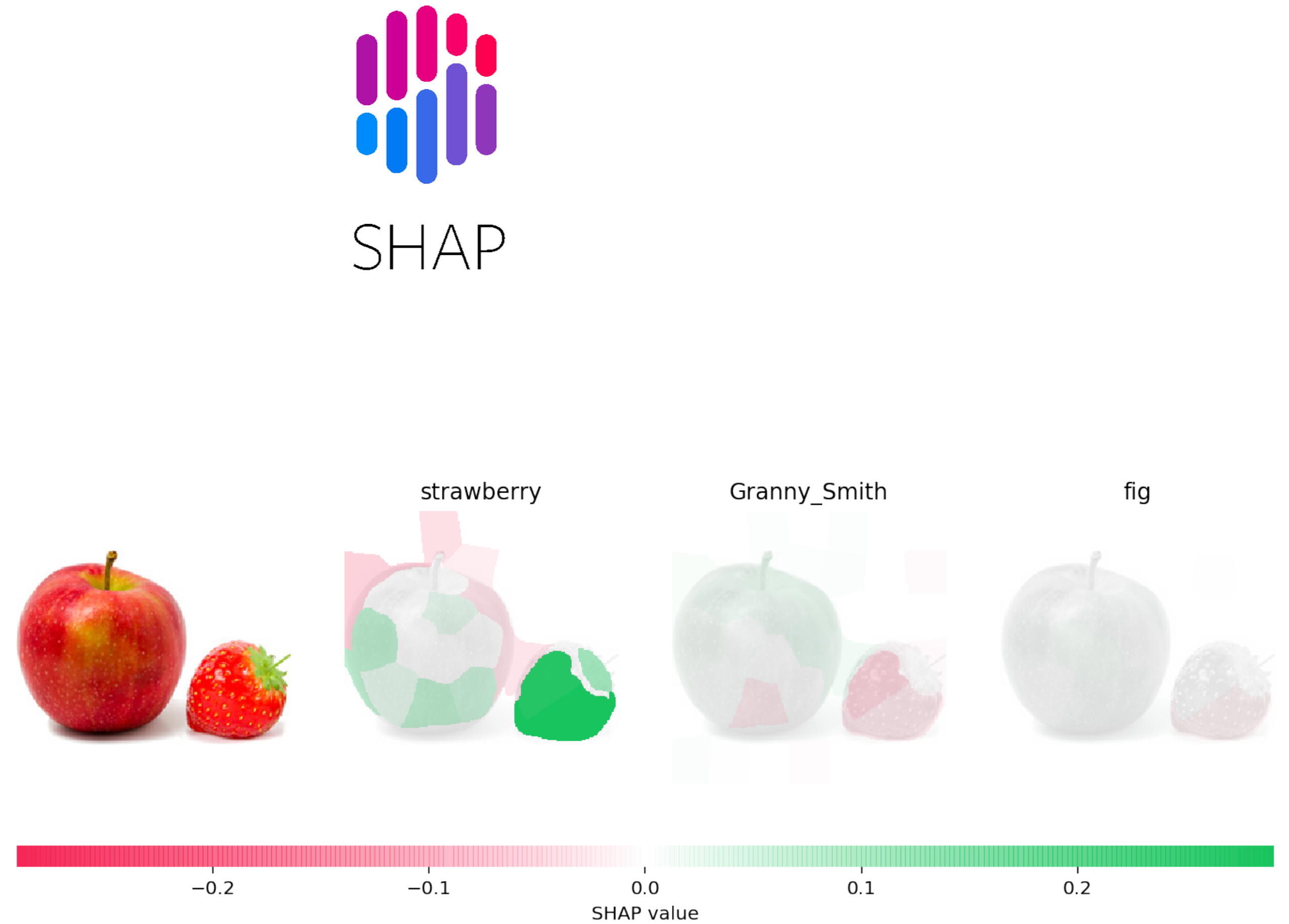
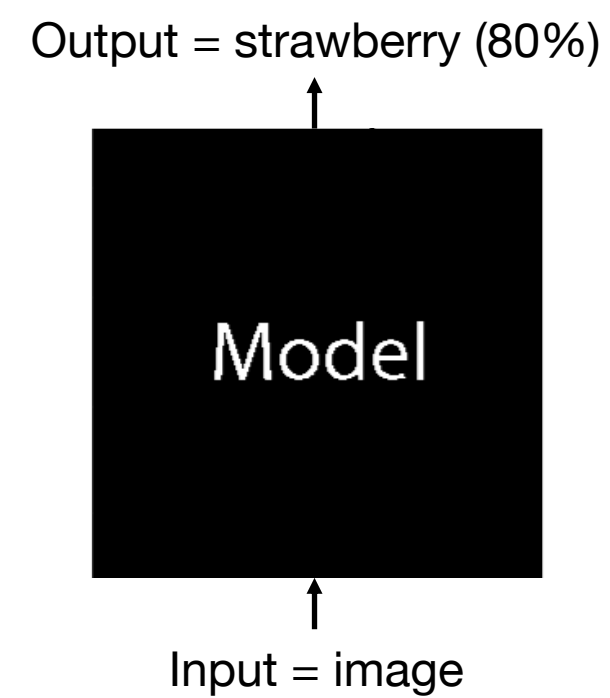
Approach: detection model

- EfficientNet as a backbone CNN
- Powerful and lightweight
- Winner's solution in Deepfake Detection Challenge



Approach: explainers

- Black-box (SHAP)



[A unified approach to interpreting model predictions, Lundberg and Lee, 2017]

Approach: explainers

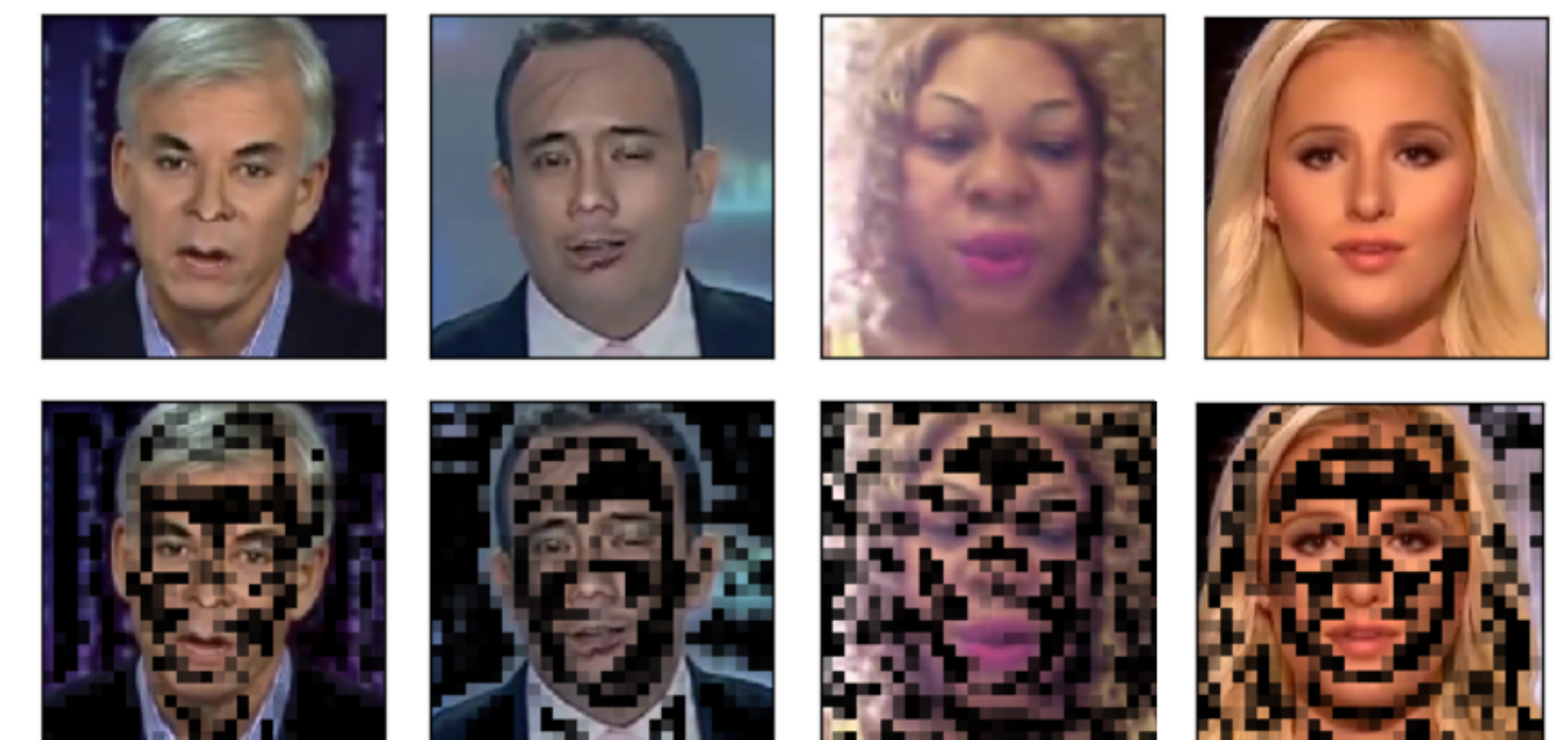
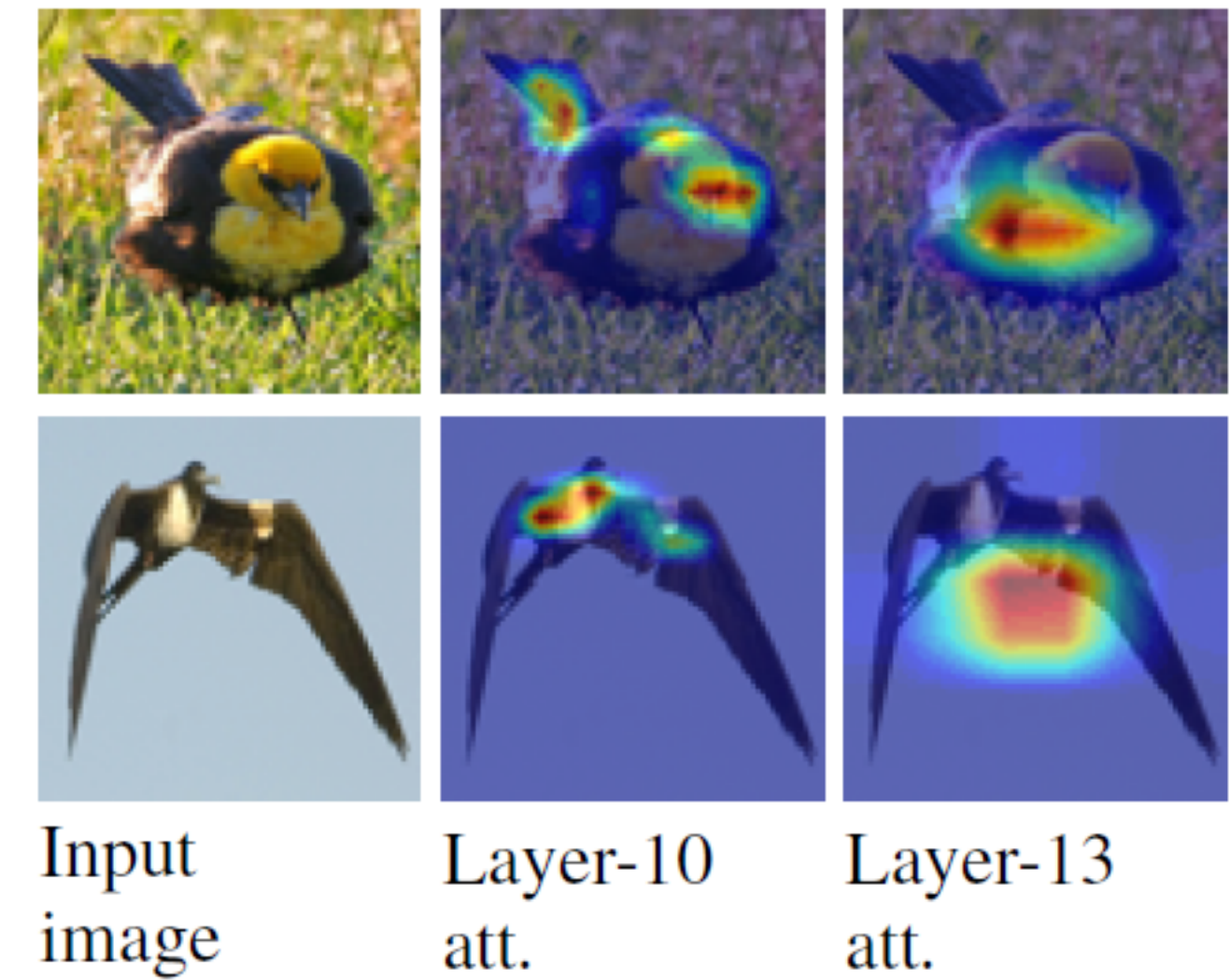
- Black-box (SHAP)
- White-box (GradCAM)



[Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization, Ramprasaath et al., 2019]

Approach: explainers

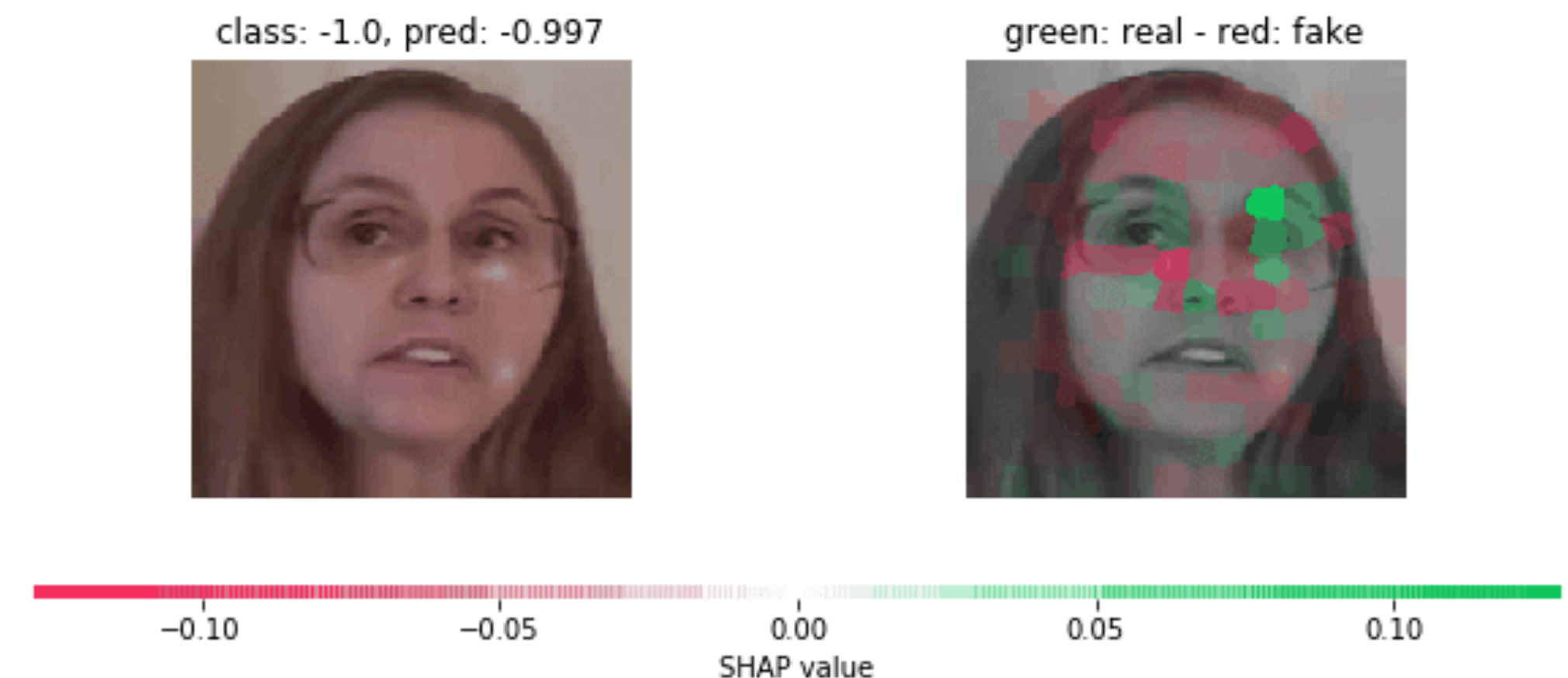
- Black-box (SHAP)
- White-box (GradCAM)
- Self-attention (LTPA, Bonettini)



[Learn to pay attention, Jetley et al., 2018]
[Video face manipulation detection through ensemble of cnns, Bonettini et al., 2020]

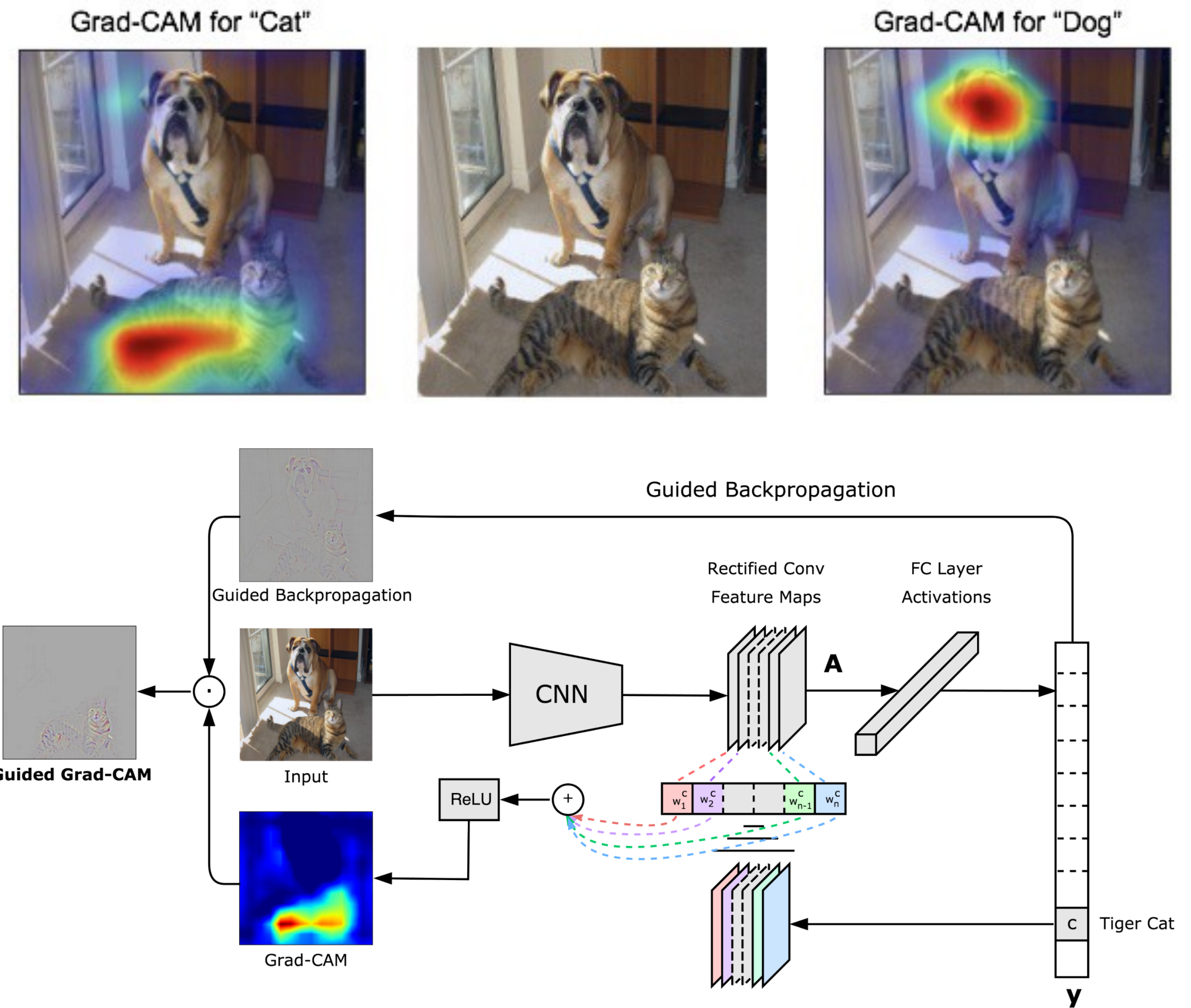
Approach: SHAP

- Model-agnostic
- Kernel SHAP for image classification
 - Segmentation
 - SHAP values assignment
- Extension: 3D segmentation for video classification



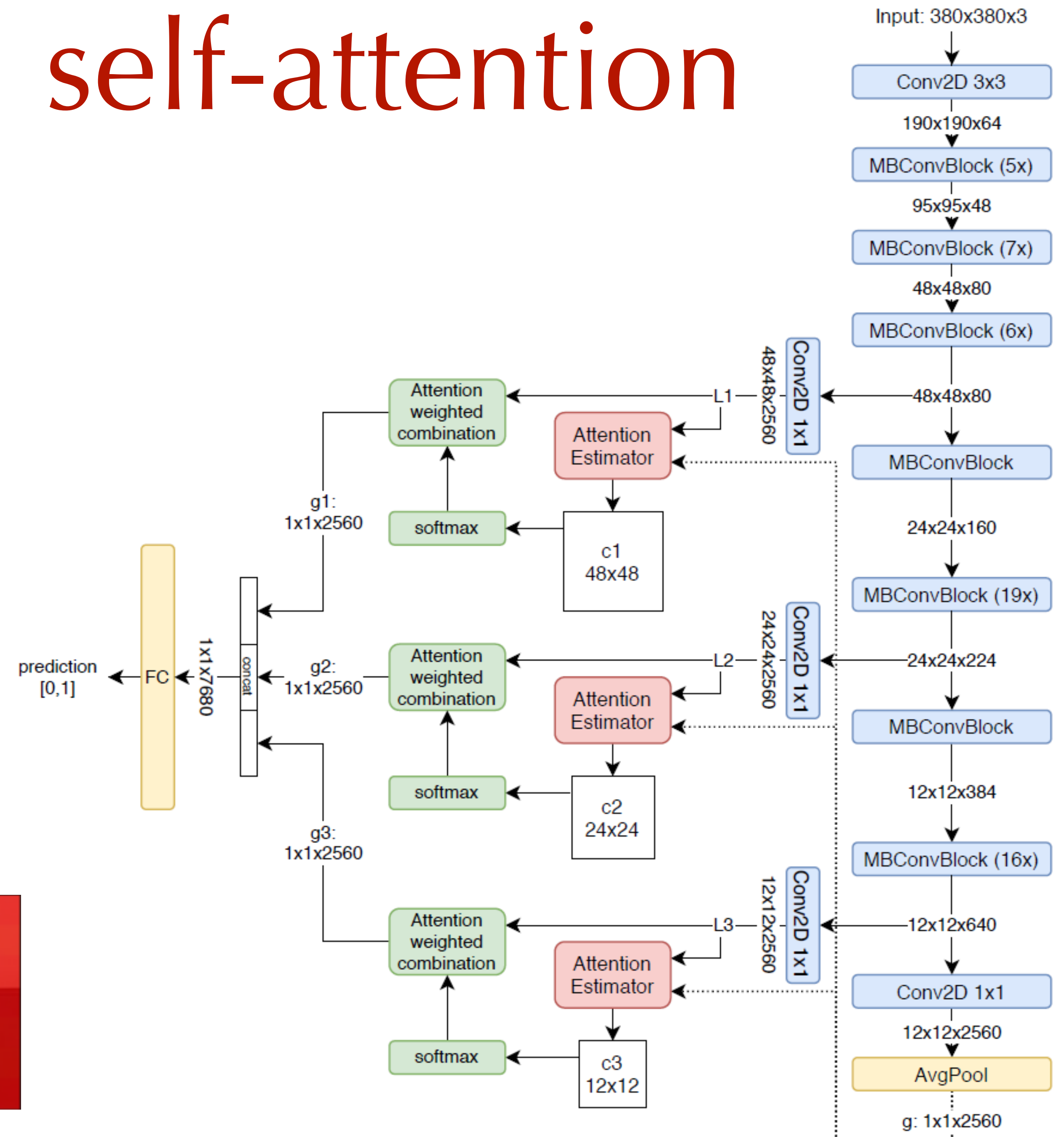
Approach: GradCAM

- Class Activation Mapping
- Neural network gradients
- Binary classification extension



Approach: self-attention

- Learn To Pay Attention (LTPA)
 - 3 attention maps



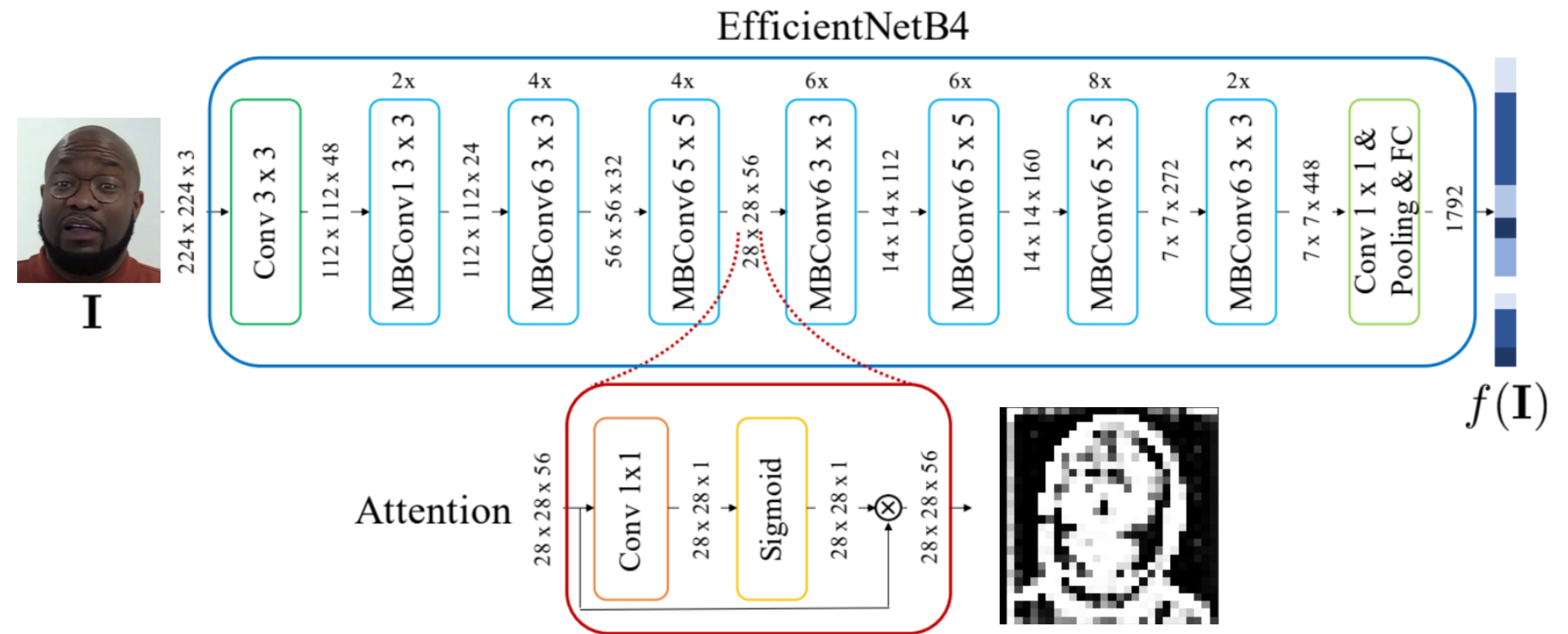
Approach: self-attention

- Learn To Pay Attention (LTPA)

- 3 attention maps

- Bonettini

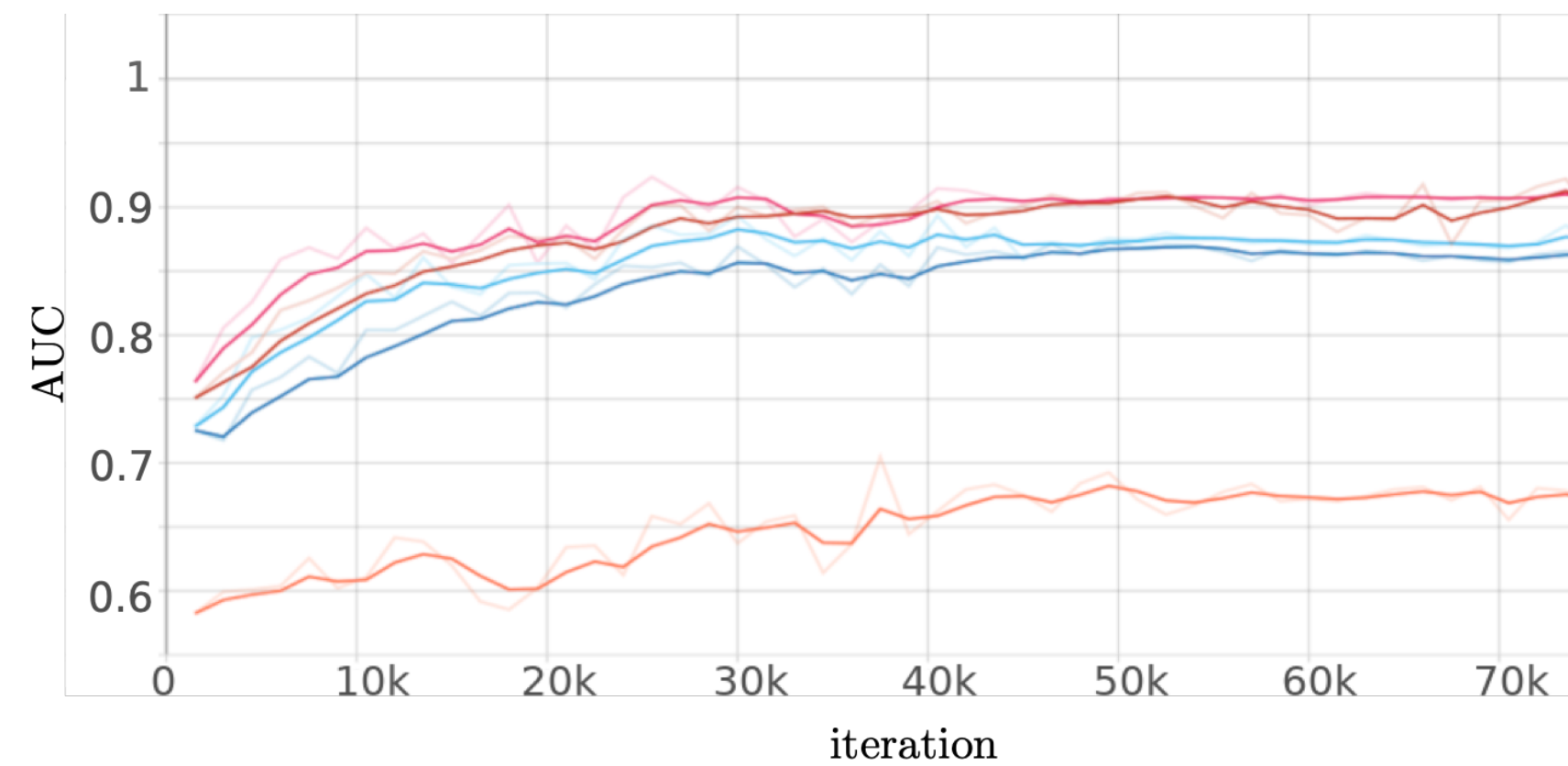
- Single attention map



Experiments

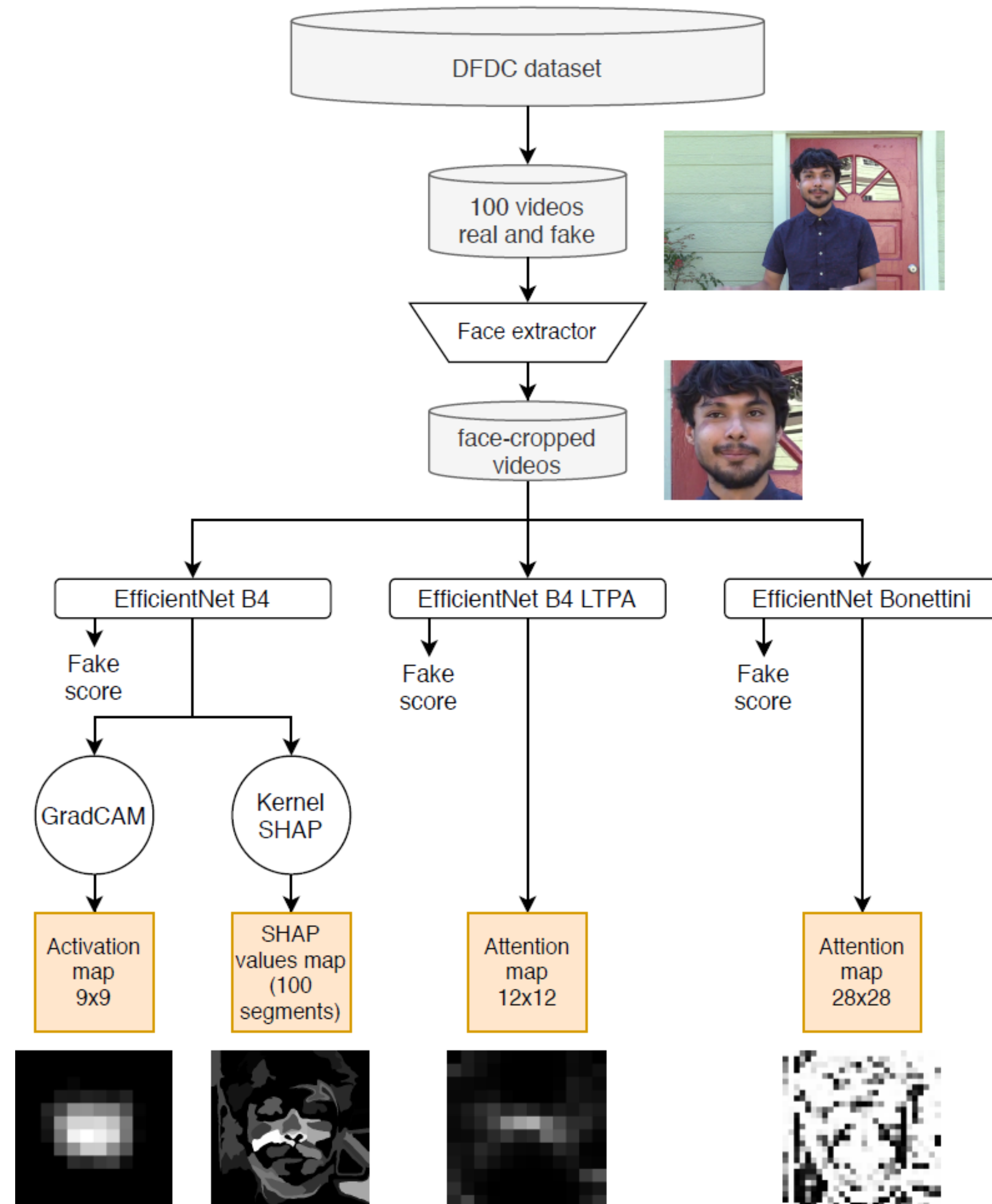
Experiments: setup

- Dataset (DFDC)
- Training

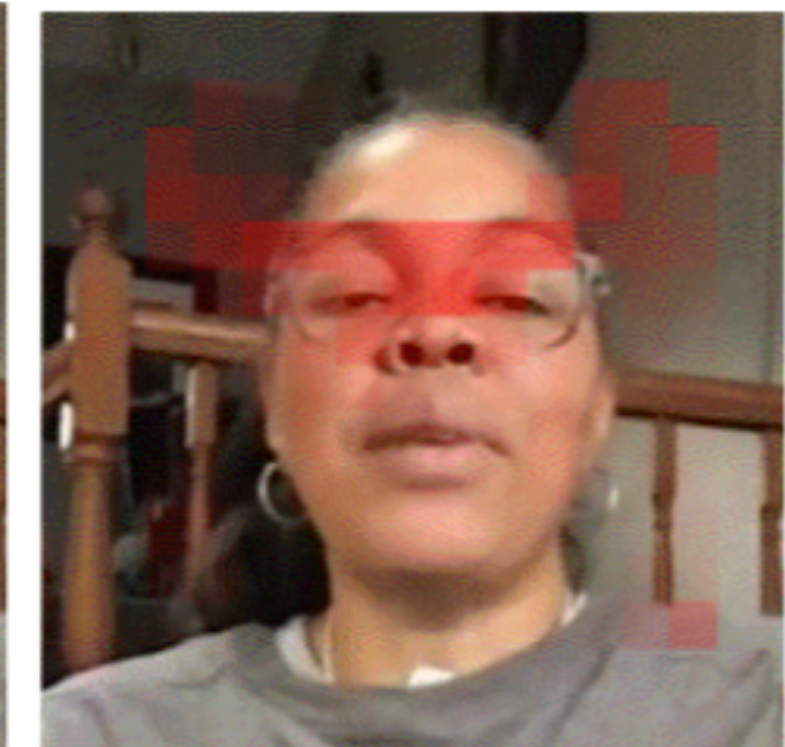


Model	Accuracy (balanced)
EfficientNet B4, 224×224	0.888
EfficientNet B4, 380×380	0.931
EfficientNet B7, 224×224	0.906
EfficientNet B7, 380×380	0.926
EfficientNet B4, LTPA, 224×224	0.879
EfficientNet B4, LTPA, 380×380	0.929
EfficientNet B7, LTPA, 224×224	0.893
EfficientNet B7, LTPA, 380×380	0.904

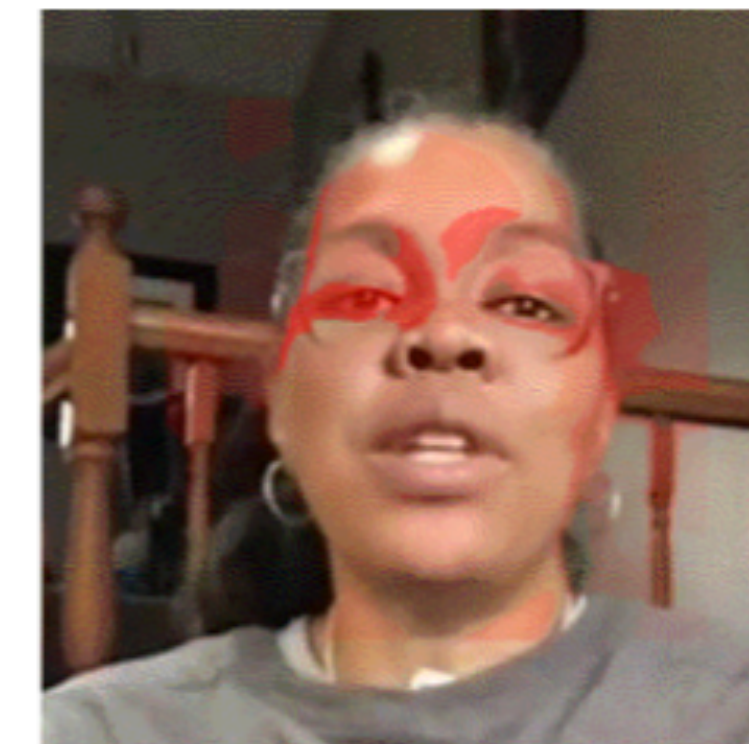
Experiments: explanations



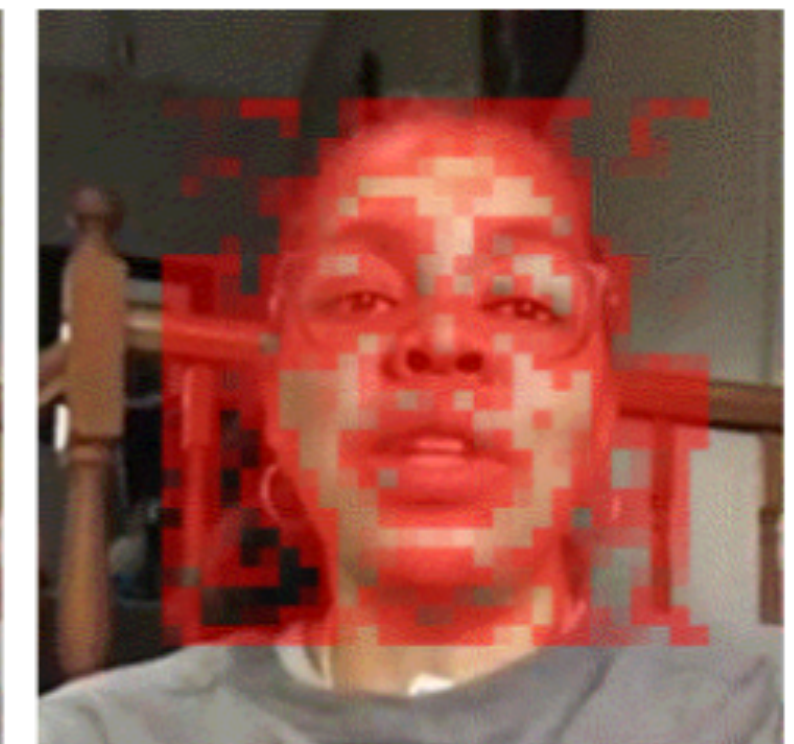
(a) GradCAM



(b) LTPA lv. 2



(c) SHAP



(d) Bonettini

Experiments: evaluation

- Metrics

- Variance

$$V = \text{avg}_{f \in \text{frames}}(\text{var}(f))$$

- Inter-frame consistency

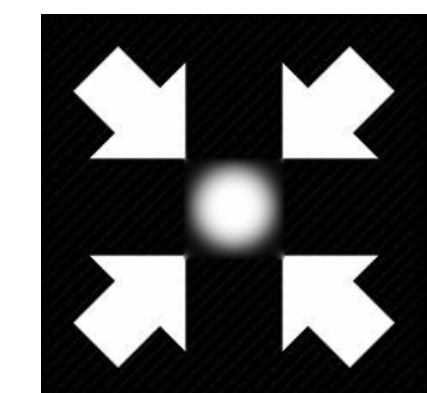
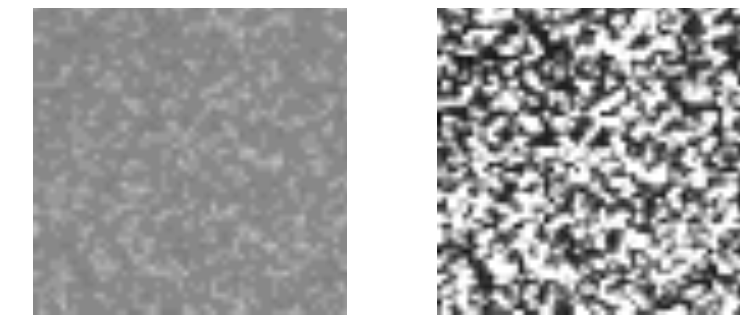
$$\tau = \text{avg}_{f \in \text{frames}}(PCC_{f,f+1})$$

- Intra-frame consistency

$$\rho = \text{avg}_{f \in \text{frames}} \left(\frac{\text{avg}_{s \in S}(a_{0.1l \cdot s}(f))}{a_{0,0}(f)} \right)$$

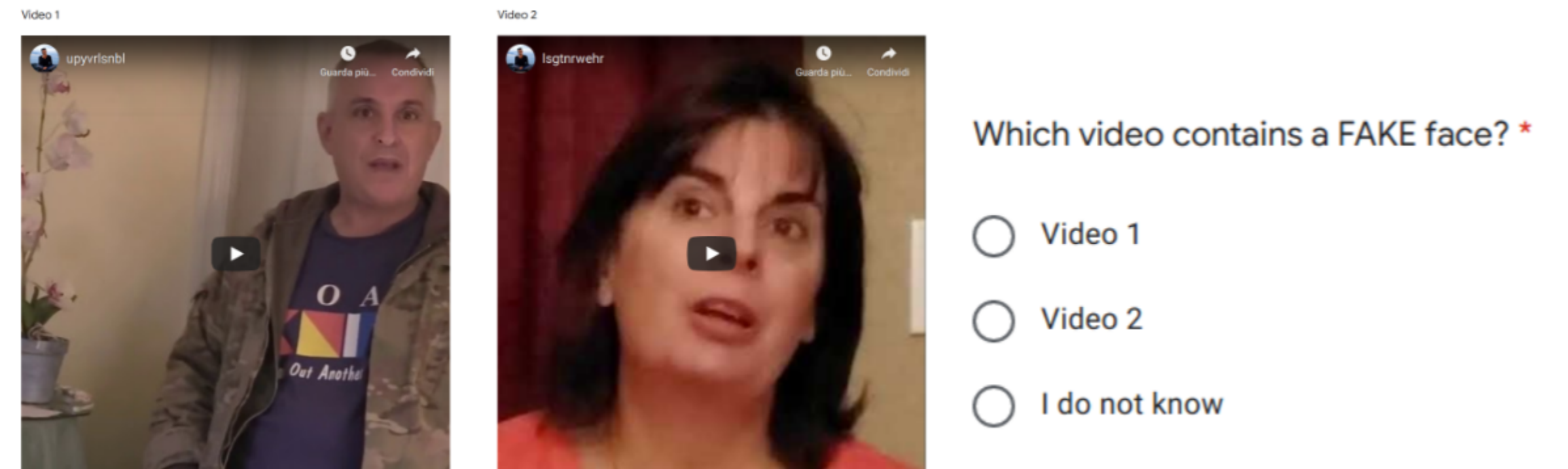
- Centredness

$$\mu = \text{avg}_{f \in \text{frames}} \left(\frac{I(\text{inner } 50\%)}{I(\text{full frame})} \right)$$



Experiments: evaluation

- User study
 - 20 real and 20 fake videos
 - 20 sections
 - 2 questions per section



Q1: A bot thinks that this face has been edited (indeed it is). In your opinion, which ones of the 4 animations best explain why the robot believes this? *



- ☐ Explanation; ☐ Explanation, ☐ Explanation. ☐ Explanation:
- ☐ Other: _____

Results


Results: metrics

- Average over 58 fake videos
- GradCAM performs best

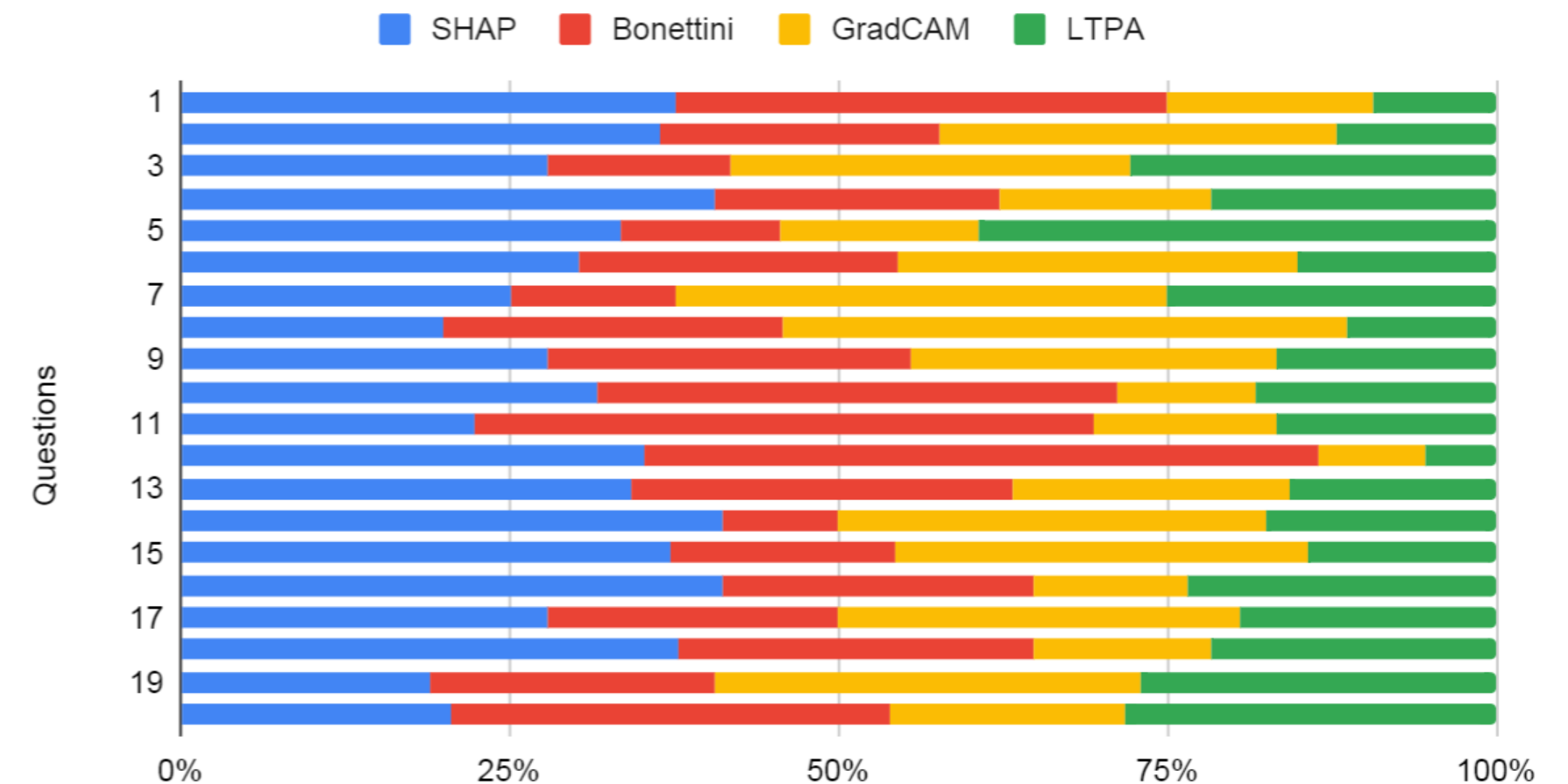


	V [0, 1]	τ [-1, 1]	ρ [-1, 1]	μ [0, 1]
Bonettini	0.0951	0.7390	0.1262	0.5286
GradCAM	0.0135	0.8756	0.7489	0.8666
LTPA	0.0108	0.7991	0.3333	0.6386
SHAP	0.0302	0.4496	0.2326	0.7348

Results: user study

- Number of answers: 67
- Accuracy: 85%
- Preferred explainer: SHAP
- Statistical “sign test” for validation
- Preference dependent on video

Number of answers		67	
Screen used	Large	43%	
	Small	57%	
Are you familiar with deepfakes?	Yes	37%	
	Heard of it	33%	
	No	30%	
Correct video identification		85%	
Explainer choices	GradCAM	165	23%
	SHAP	221	31%
	LTPA	137	19%
	Bonettini	185	26%



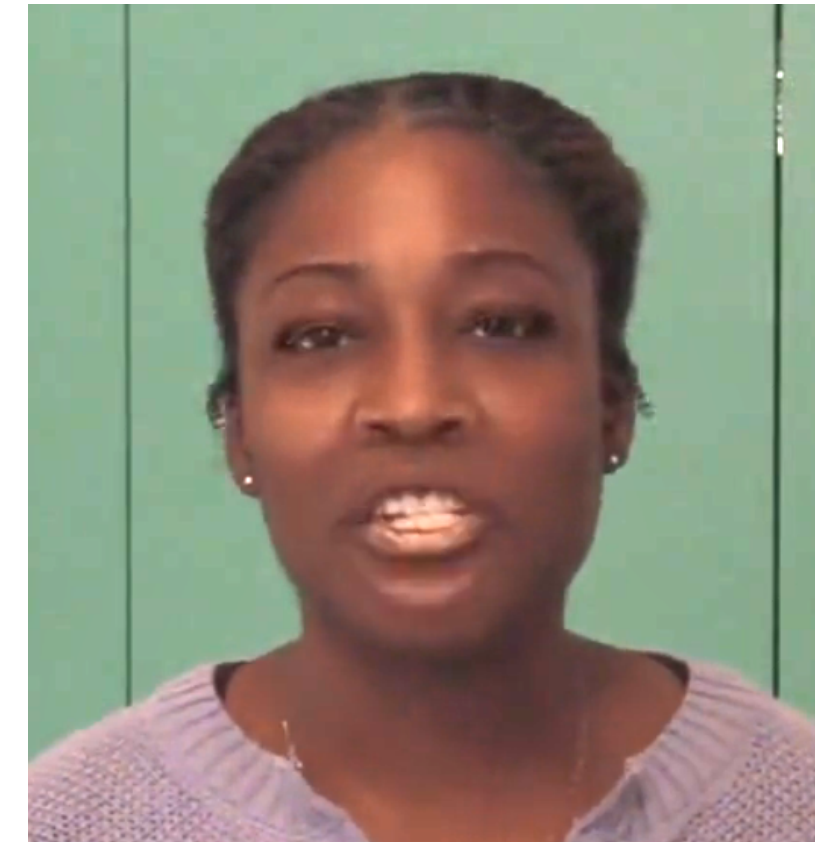
Conclusion

Conclusion

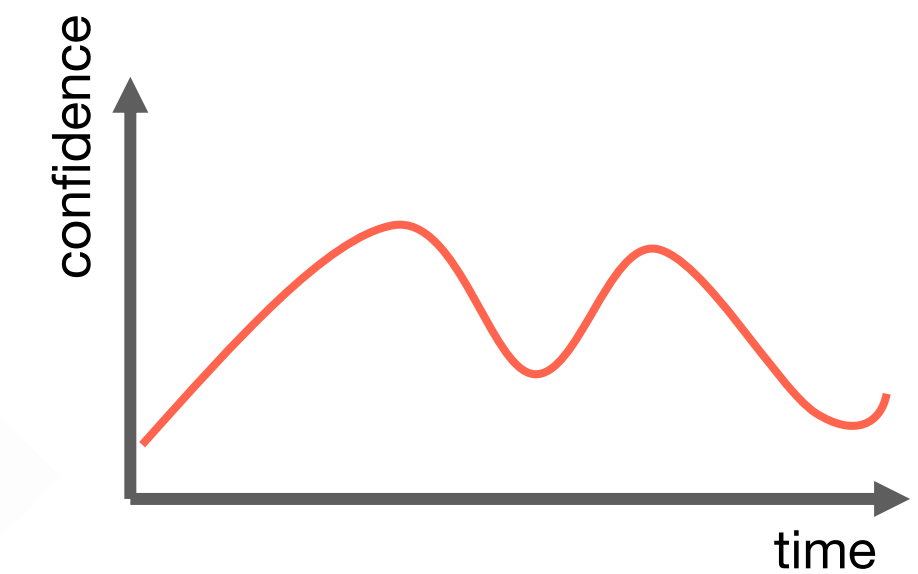
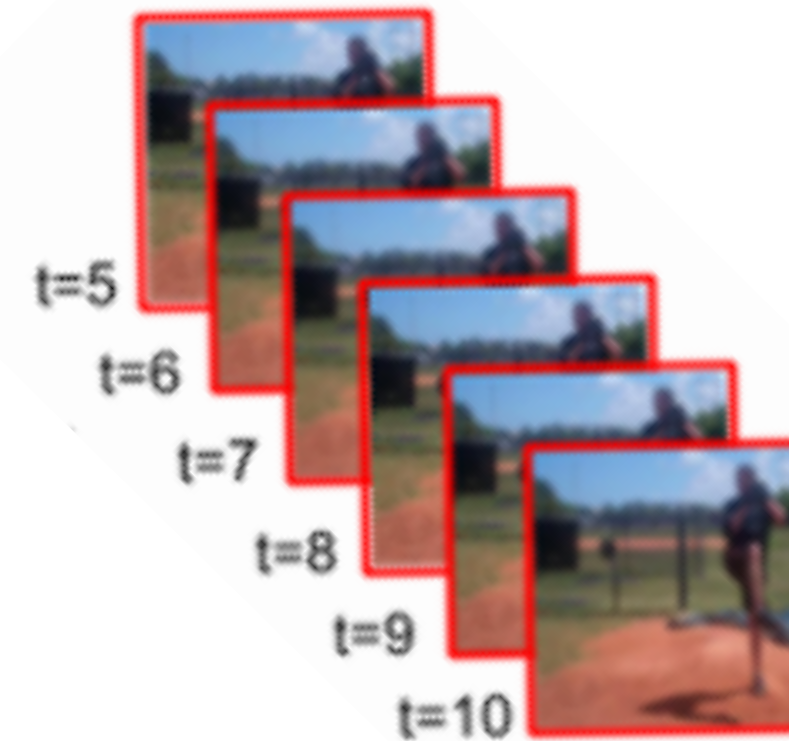
- We implemented and extended 4 explanation techniques
- We defined intrinsic and extrinsic metrics
- We empirically compared the explainers based on them
- We performed a user survey
- Human perception is not always aligned with objective metrics

Future work

- Captioning explanation maps
- Weighting explanations on classifier's confidence



« LIPS ARE FAKE! »



Thank you