

# Capsule-Forensics: Using Capsule Networks to Detect Forged Images and Videos

Huy H. Nguyen (SOKENDAI, Japan)

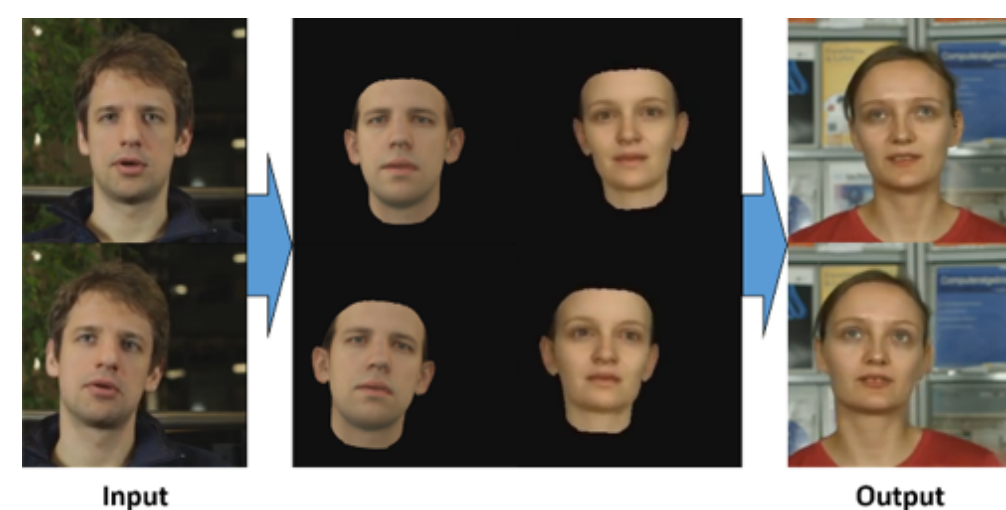
Junichi Yamagishi (NII, Japan)

Isao Echizen (NII, Japan)

## Generating of Fake Videos Impersonating a Person Using Deep Learning



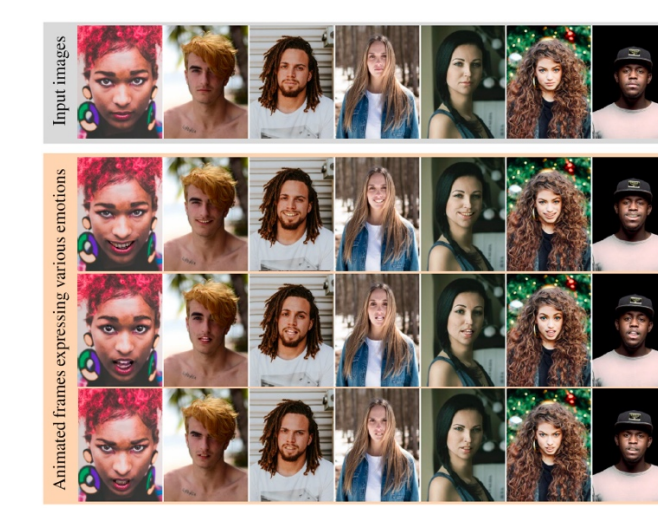
Face2Face: Real-time facial reenactment (Thies et al. 2016)



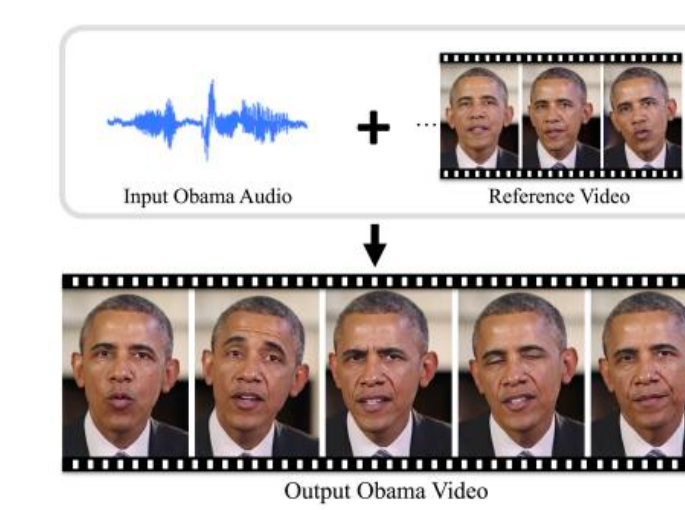
Deep Video Portraits: Face2Face + Translating head poses (Kim et al. 2018)



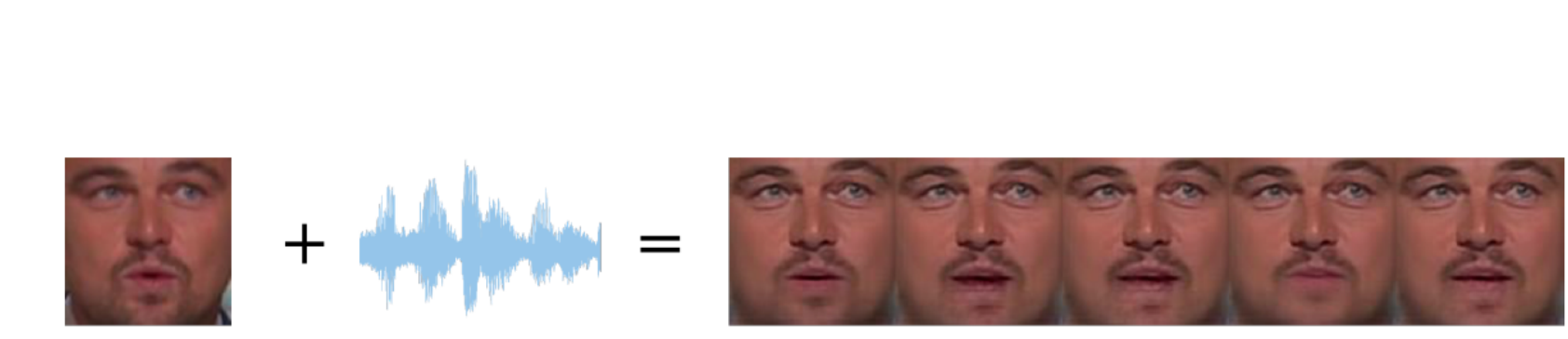
Deepfakes Video face swapping (2017)



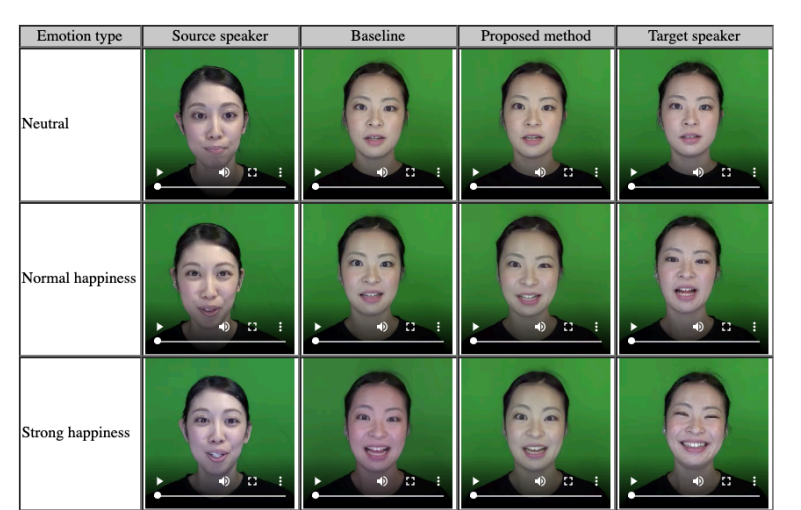
Bringing portraits to life (Averbuch-Elor et al. 2017)



Synthesizing Obama: Learning lip sync from audio (Suwajanakorn et al. 2017)



Speech2Vid (Chung et al. 2017)



Audiovisual speaker conversion (Fang et al. ICASSP 2019)

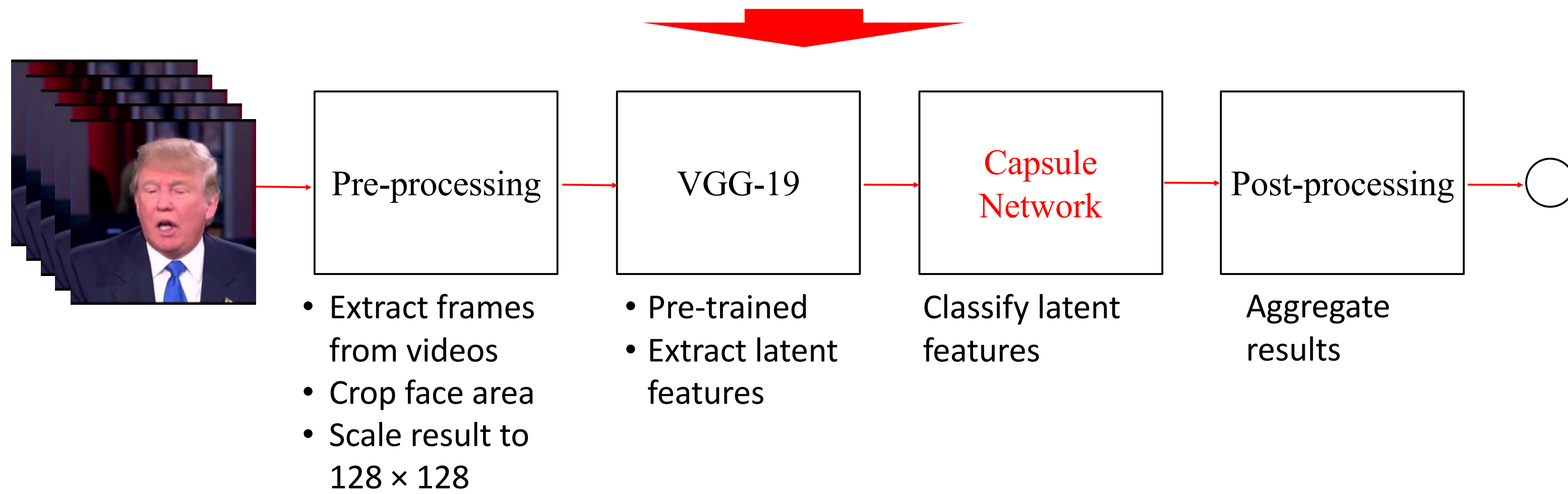
➤ Deep learning methods enable **non-expert users** with **ordinary PCs** to create **realistic** forged images and videos by using **data available on social networks**.

➤ Materials required for generating fake videos **have been simplified** over time.

➤ Forgery detectors need to be **regularly updated** to deal with

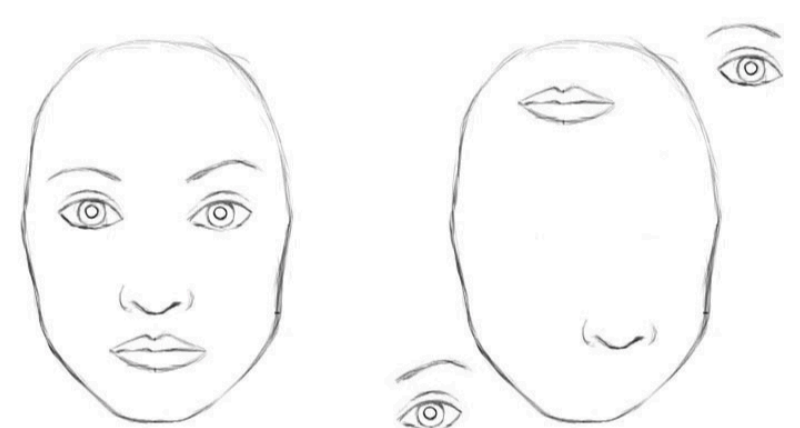
- New kind of attacks
- Better quality forged images/videos

➔ Is there a general framework that could be applied for any kind of attack???



## CNN vs. Capsule Network

• In **computer vision perspective**, convolutional neural networks (CNNs) has **viewpoint invariant** property but **lacking of** information about **relative spatial relationships** between features → Capsule network can solve this problem.

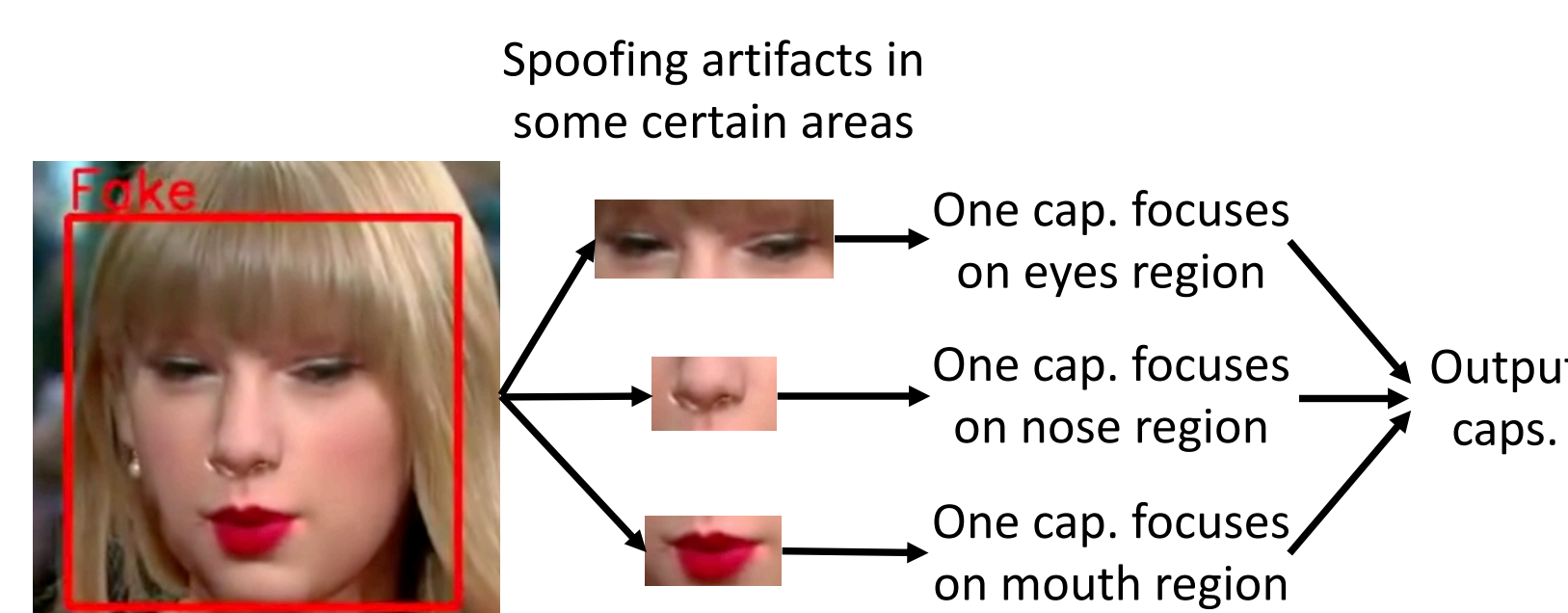


The two pictures are similar in the perspective of a CNN but dissimilar in the view of a capsule network.  
Source: Max Pechyonkin.

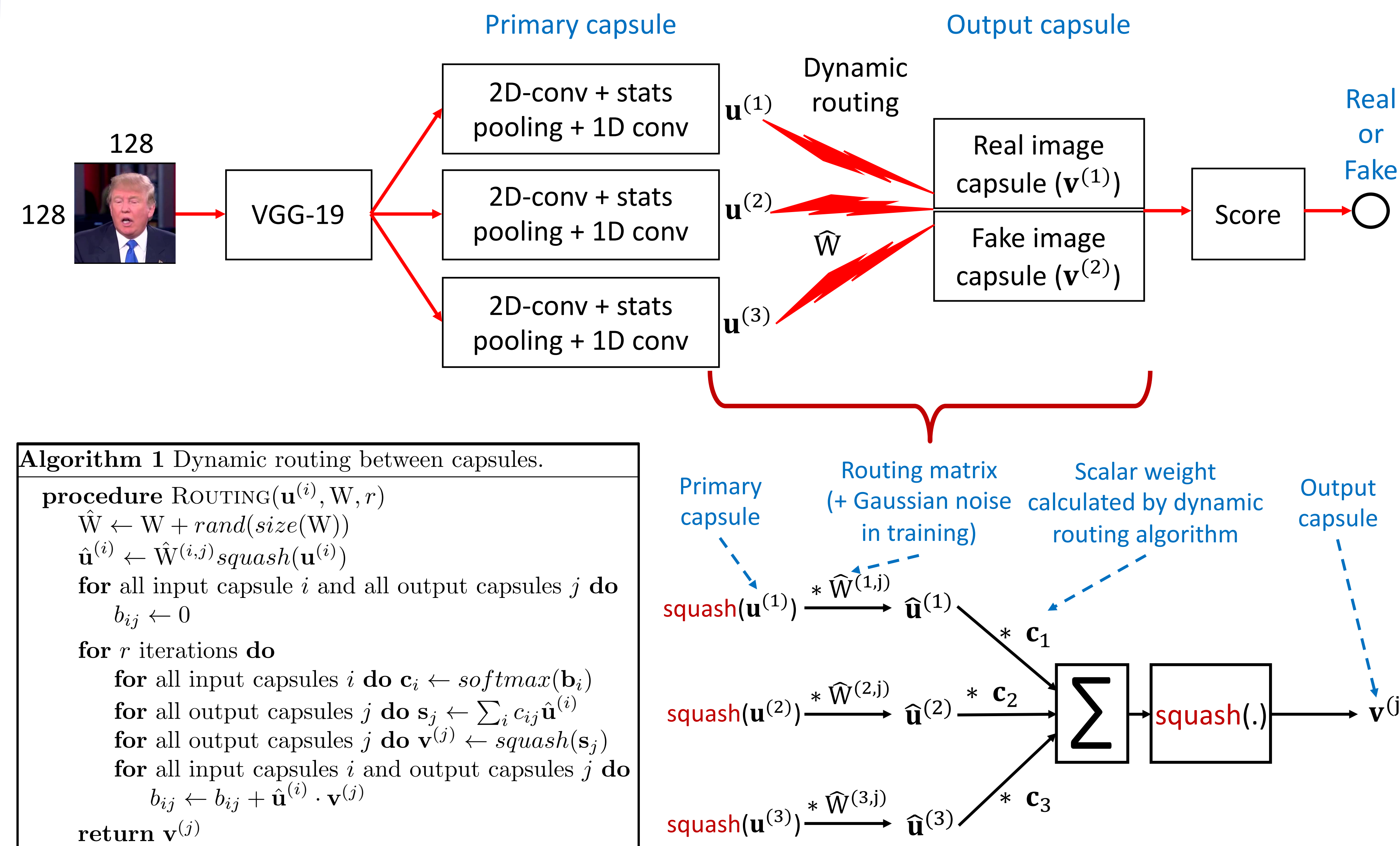
• Capsule networks have several capsules, each capsule is a **CNN** learning some **specific** representations.

• The **agreements** between low-level capsules decides the **activations** of the high-level capsules.

• In **digital image forensics perspective**, each low-level capsule may capture some specific representations of **spoofing artifacts** in some certain area, or some specific kind of **irregular noises** created by presentation attacks.



## Capsule Network

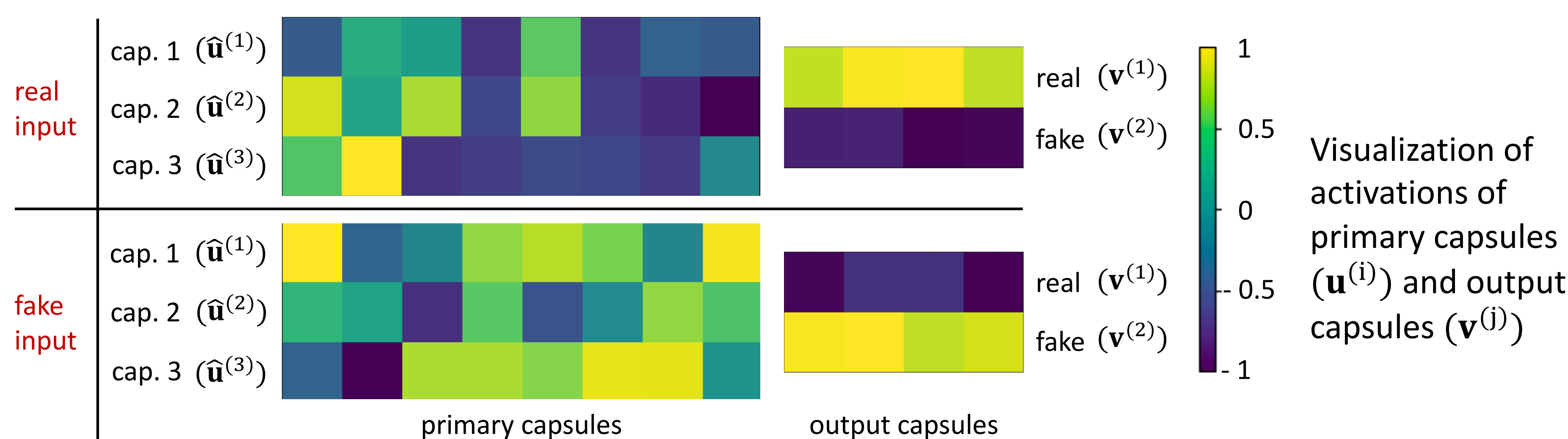


Squash function, used to scale vector magnitude to unit length:

$$\text{squash}(\mathbf{u}) = \frac{\|\mathbf{u}\|_2^2}{1 + \|\mathbf{u}\|_2^2} \frac{1}{\|\mathbf{u}\|_2} \mathbf{u}$$

Score function, used to determine the predicted label probabilities:

$$\hat{y} = \frac{1}{m} \sum_i \text{softmax} \left( \begin{bmatrix} \mathbf{v}^{(1)\top} \\ \mathbf{v}^{(2)\top} \end{bmatrix}_{:,i} \right)$$



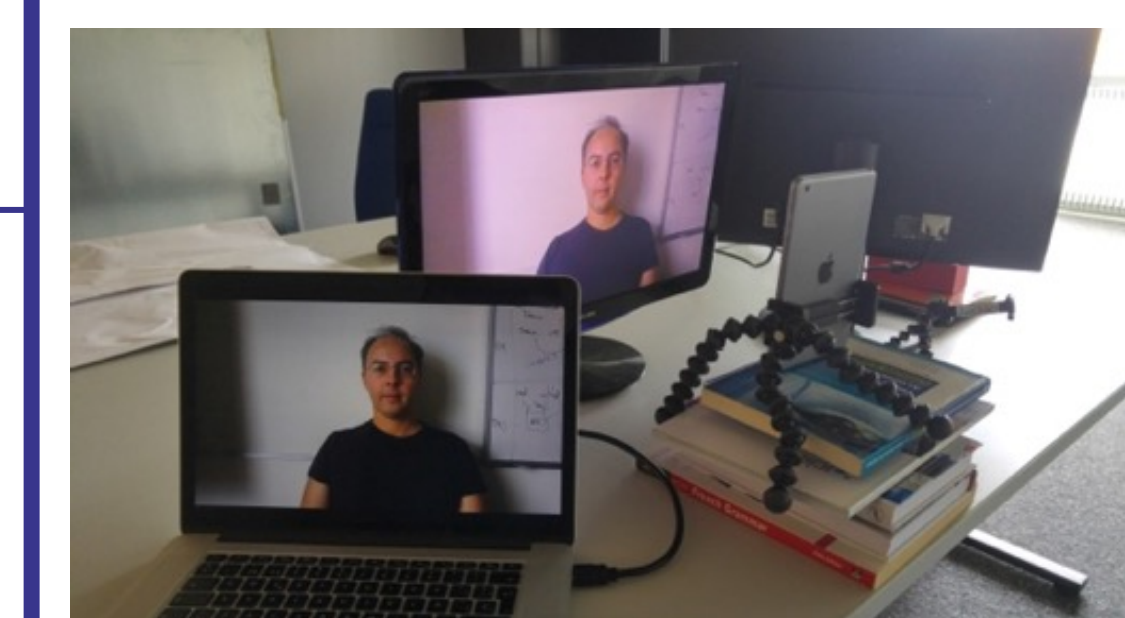
## Evaluation

Type of Attack	Detection Accuracy (%)
Replay Attack	100.00 *
CG vs. Natural Images	100.00 *
Deepfakes (Frames)	95.93 *
Deepfakes (Video)	99.23 *
Face2Face (c0 - Frames)	99.37
Face2Face (c0 - Video)	99.33
Face2Face (c23 - Frames)	96.50
Face2Face (c23 - Video)	96.00
Face2Face (c40 - Frames)	81.00
Face2Face (c40 - Video)	83.33

**Note:**

- c0, c23, c40: H264 compression levels
- number \*: state-of-the-art result

Some examples from the evaluation datasets:



Idiap Replay Attack database (Chingovska et al. 2012)



Facial reenactment (Face2Face) in the FaceForensics database (Rössler et al. 2018)

