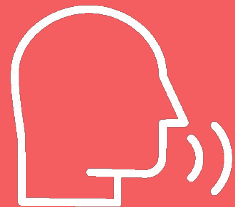


# **Audio & Image classification, Image retrieval**

Progetto finale  
Digital Signal & Image Management

# Audio classification



“Giacomo”

“Leonardo”

“Riccardo”

# Dataset

Task: identificare chi sta parlando, tra i componenti del gruppo  
L'interlocutore deve pronunciare 3 cifre (senza altri vincoli)

Costruzione dataset:



train set di 100 osservazioni (+ 20 di validation);      test set di 33 osservazioni

# Modelli

1. NN senza features pre-calcolate
2. NN con features audio base
3. NN con features da modello pre-trainato

↓ complessità  
features  
crescente

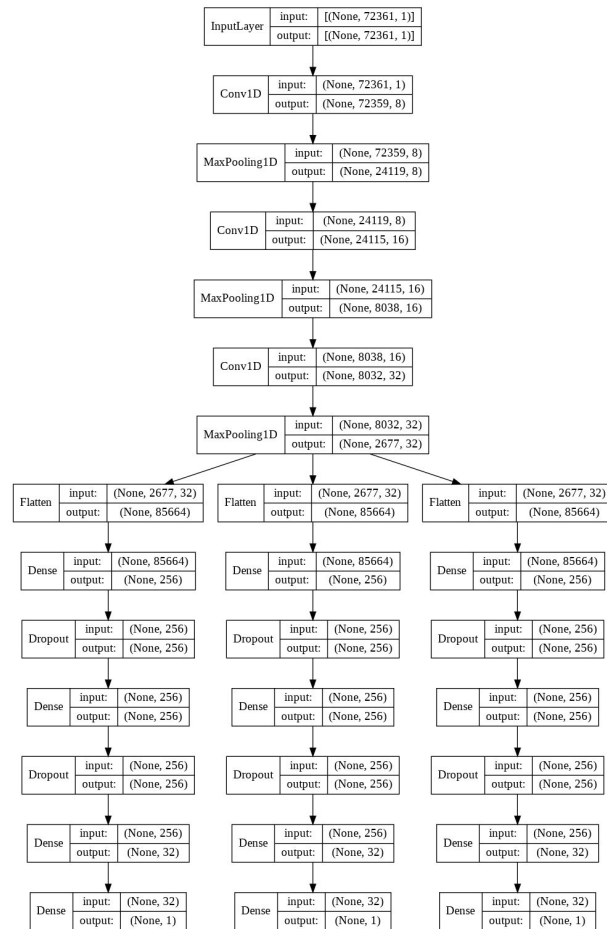
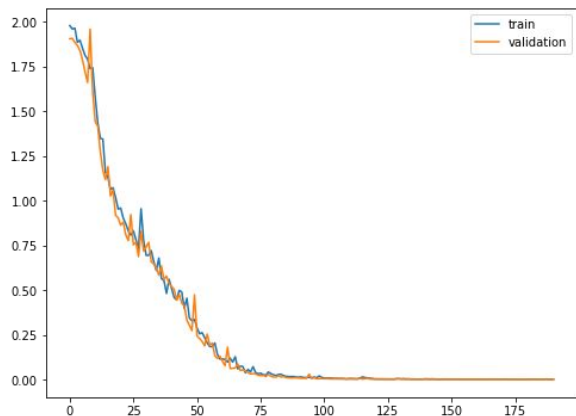
Per tutti i modelli, sono state scelte 3 sigmoidi, rispetto ad 1 softmax

- si riesce a discriminare meglio quando a parlare non è nessuno tra i componenti
- si ottengono risultati migliori

# Modello 1 (NN senza features pre-calcolate)

2 fasi:

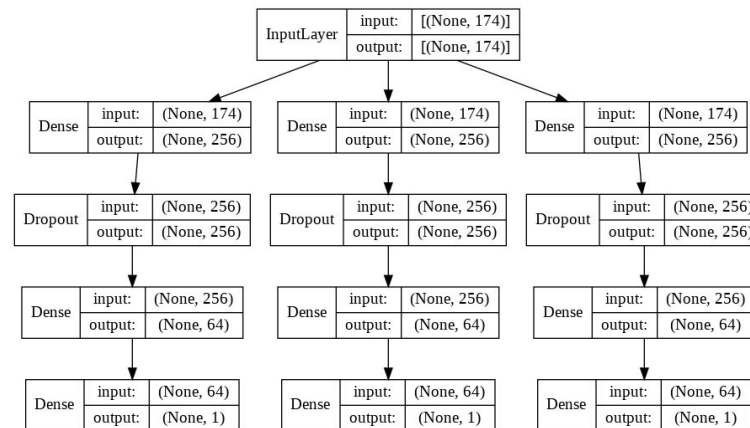
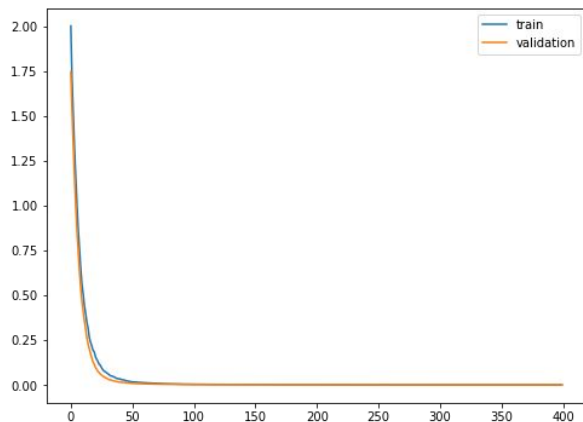
- Convoluzioni e MaxPooling (a 1 dimensione)
- Flatten, Dense e Dropout



# Modello 2 (NN con features audio base)

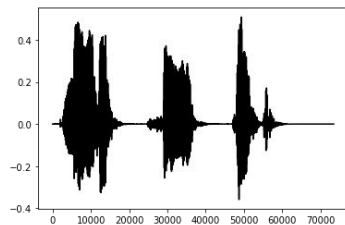
Librosa features → utilizzate quelle che hanno dimensione indipendente da lunghezza audio

- [mfcc, chroma\_stft, melspectrogram, spectral\_contrast, tonnetz, energy]

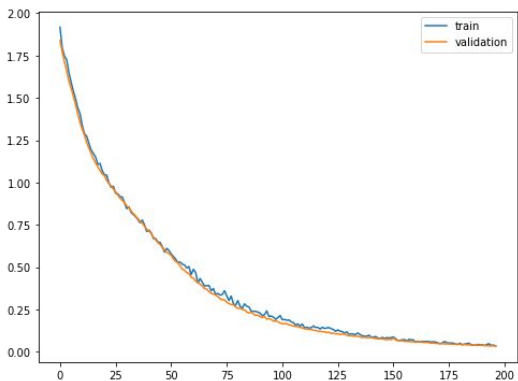
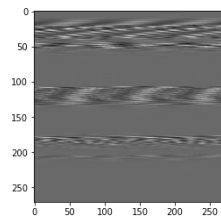


# Modello 3 (NN con features da modello pre-trainedo)

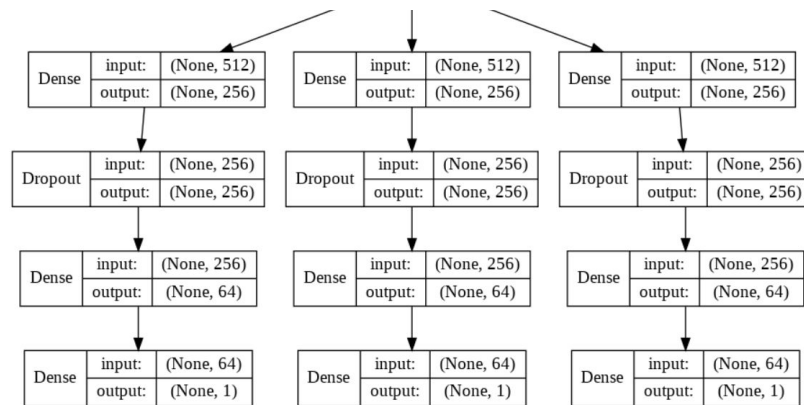
segnale 1D



reshape 2D e resize



VGG16 freezed model



# Risultati

F1-score	Modello 1	Modello 2	Modello 3
Giacomo	0.17	0.91	0.85
Leonardo	0.51	0.95	0.59
Riccardo	0.53	0.87	0.78
	<b>0.40</b>	<b>0.91</b>	<b>0.74</b>

Tempo di training	480s	130s	420s
-------------------	------	------	------



# Image classification



“Giacomo”

“Leonardo”

“Riccardo”

# Workflow

Task: identificare il volto di una persona , tra i componenti del gruppo.  
Assegnare probabilità predizione classificatore.

## Costruzione Dataset:

Scattate foto per ogni  
componente da diverse  
posizioni e angolazioni



Utilizzo di un Face  
Detector per ritagliare  
i volti (pulizia  
supervisionata)



Partizionamento in  
train, validation e test

Modello reti neurali  
MobileNetv2

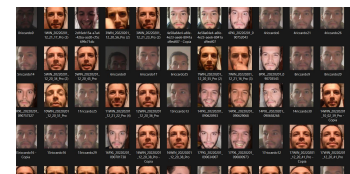
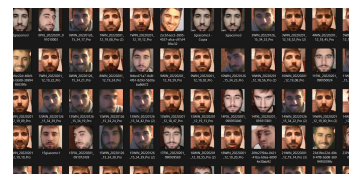
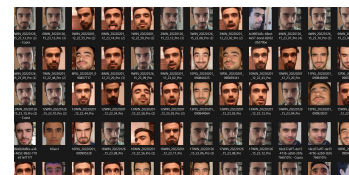
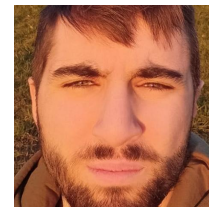
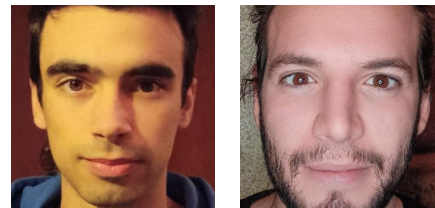
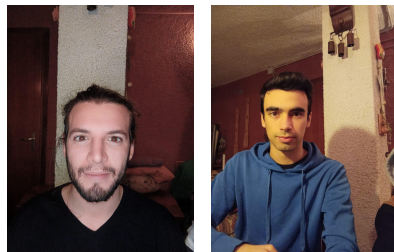


Modello reti neurali  
VGG16



Modello LBPH (Local Binary  
Pattern Histogram)

# Dataset

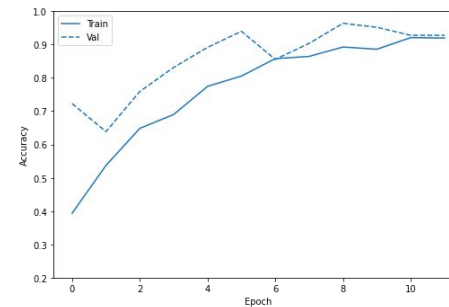
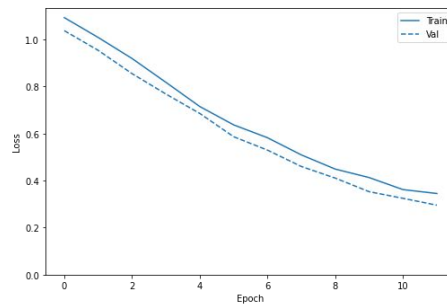
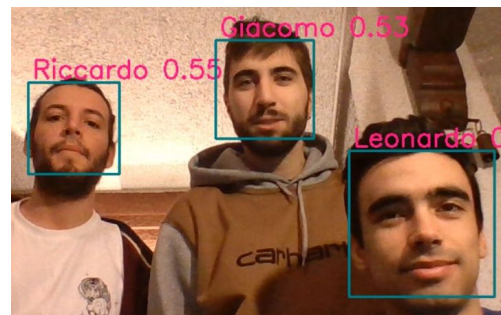
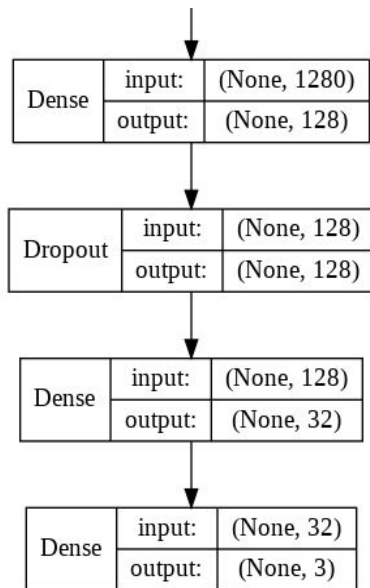


Train: 672 images

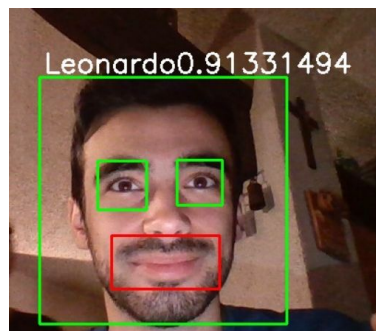
Validation: 73 images

Test: 83 images

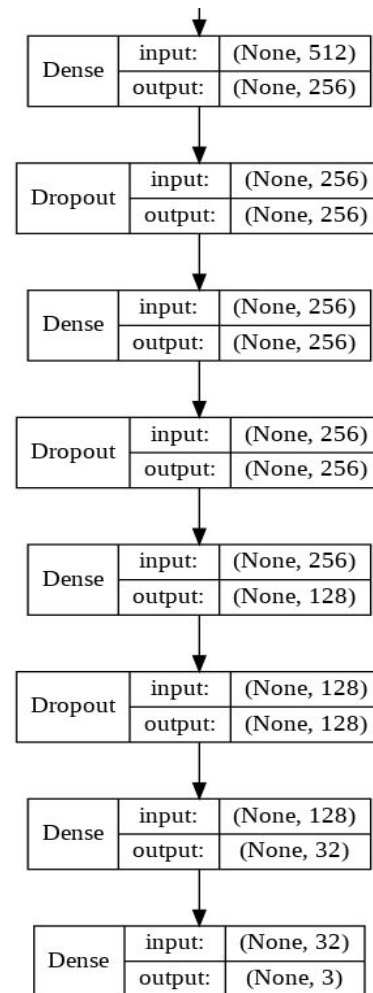
# Modello 1 (MobileNetV2)



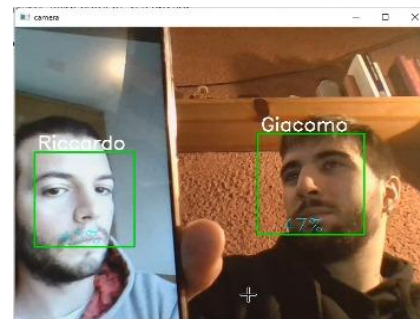
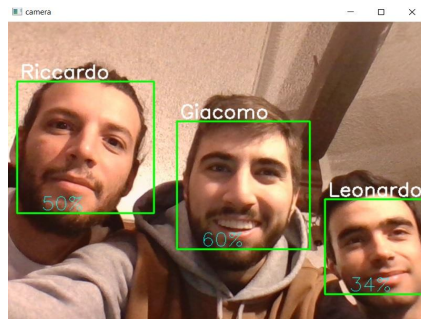
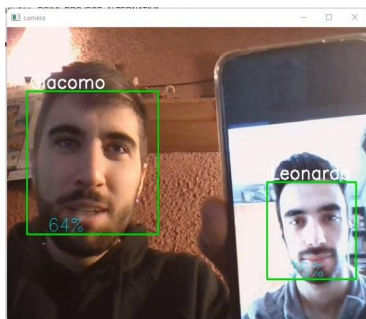
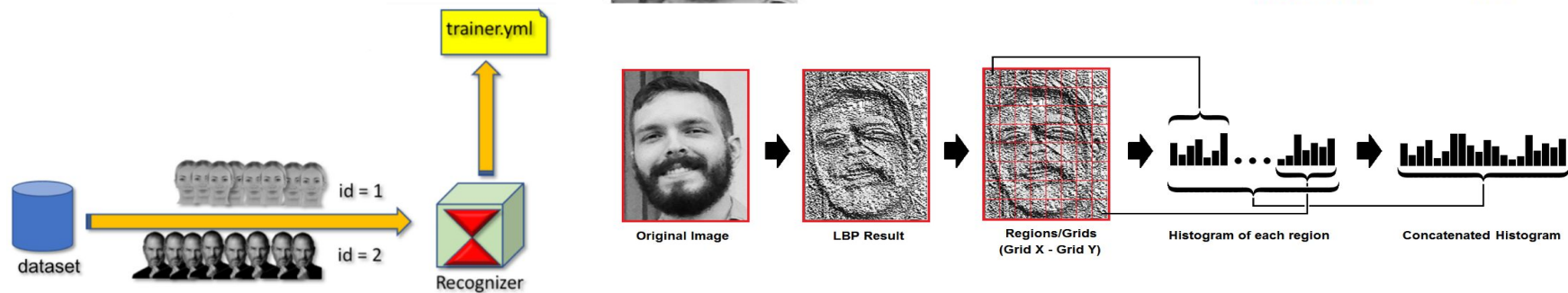
# Modello 2 (vgg16)



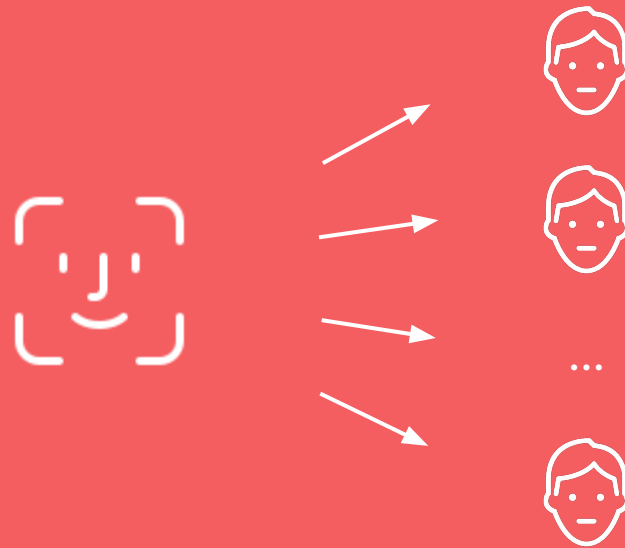
	precision	recall	f1-score	support
0	0.92	1.00	0.96	23
1	1.00	0.94	0.97	32
2	0.96	0.96	0.96	28
accuracy			0.96	83
macro avg	0.96	0.97	0.96	83
weighted avg	0.97	0.96	0.96	83



# Modello 3 (LBPH)



# Image retrieval

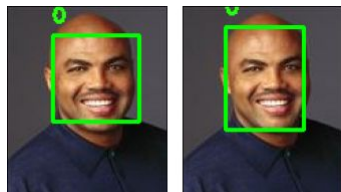


# Dataset

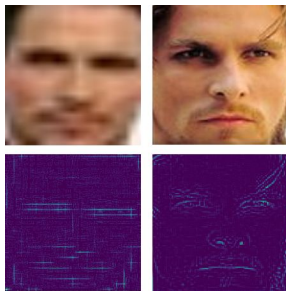
Task: trovare le 10 facce più simili, tra quelle presenti nel dataset, all'immagine in input

Costruzione dataset:

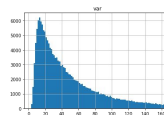
face detection  
tramite  
MTCNN



rimozione  
immagini low  
quality



estrazione  
feature

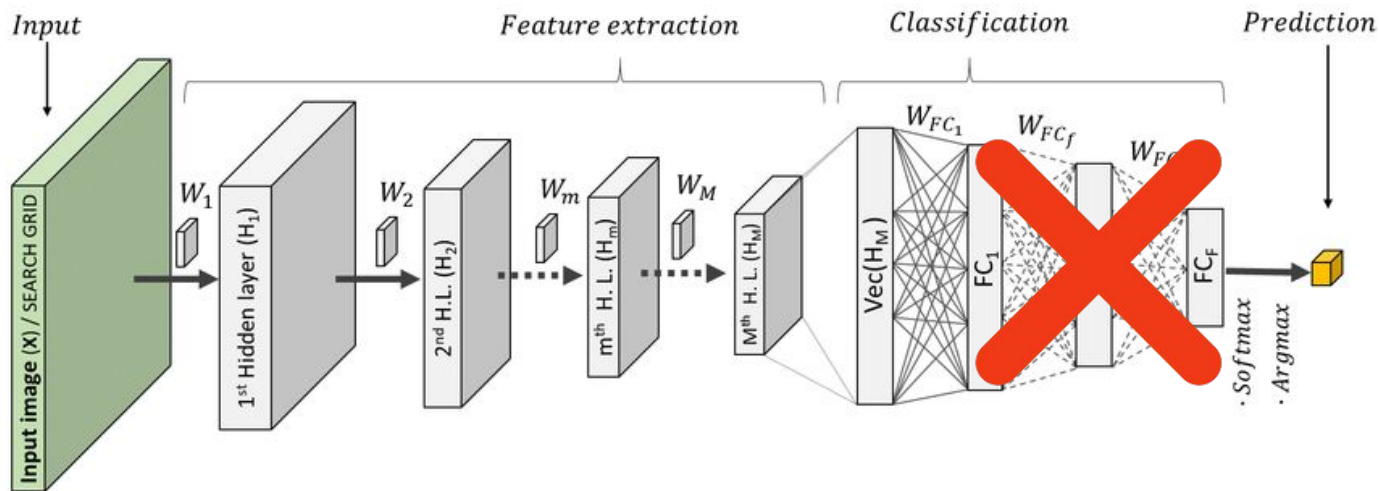


$[[x_1, x_2, x_3, x_4, x_5, \dots]]$

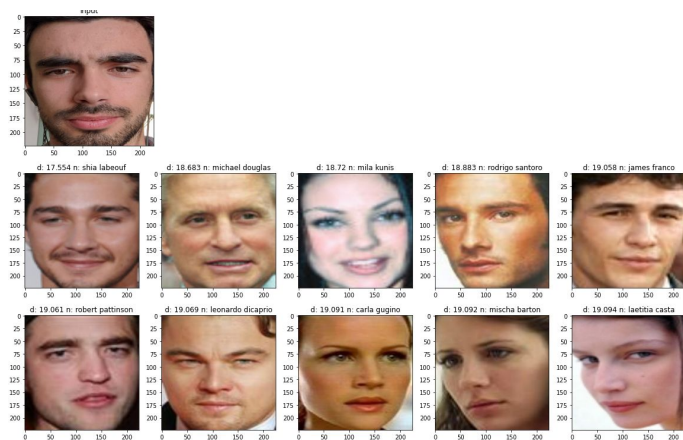
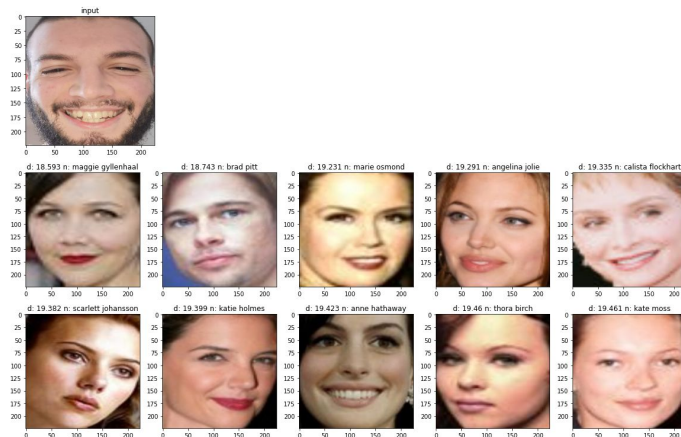
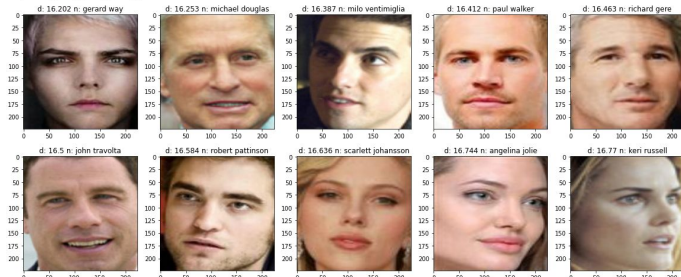
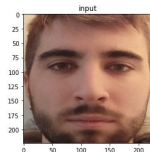


# Modelli

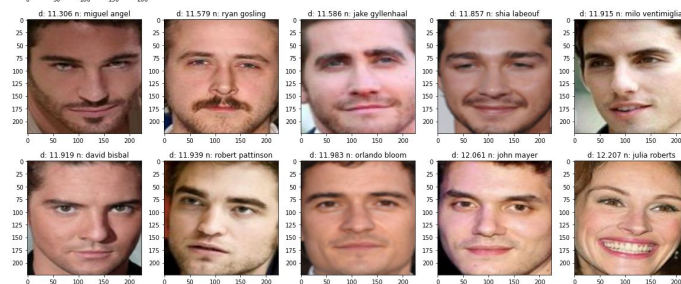
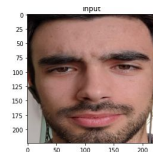
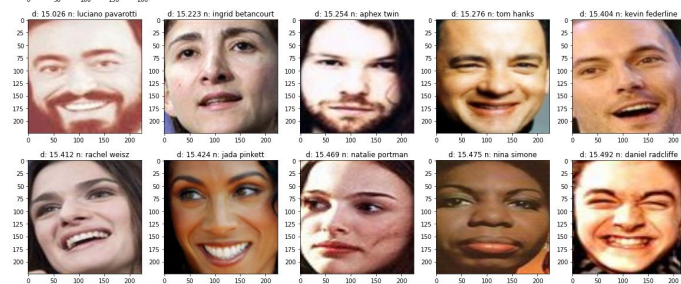
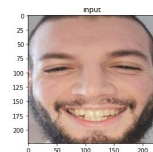
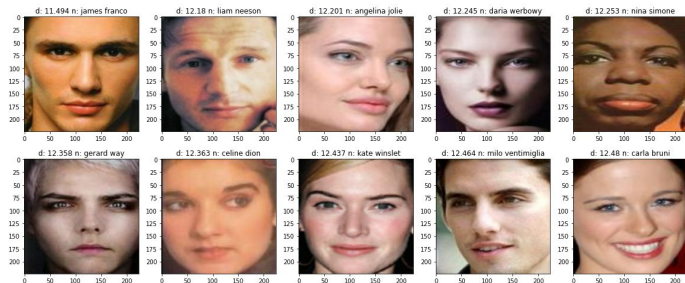
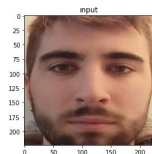
- Mobilenet V2
  - DenseNet
  - VGGFACE2 con ResNet50
- *pesi Imagenet*
- *pesi VggFace*



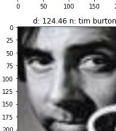
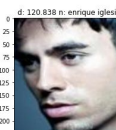
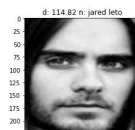
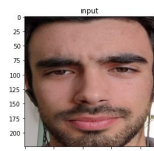
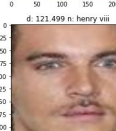
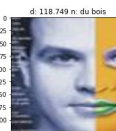
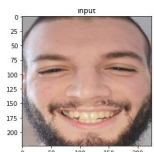
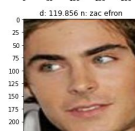
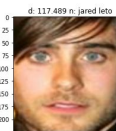
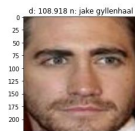
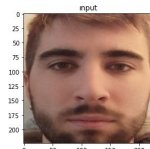
# Modello 1 (MobileNetV2)



# Modello 2 (DenseNet)



# Modello 3 (VggFace2)



A large red square with a white border. Inside the square, the words "THE" and "END" are written in white, bold, sans-serif capital letters, stacked vertically in the center.

**THE  
END**