# Interpreting CNN Models for Apparent Personality Trait Regression

#### Carles Ventura, David Masip, Agata Lapedriza



SUnA SCENE UNDERSTANDING & ARTIFICIAL INTELLIGENCE Computer Vision & A.I Recognition for Objects and Gestures

# Outline

- Introduction
- Related Work
- Experiments
  - Images + audio vs Images for personality trait regression
  - Finding Discriminative Regions in video frames
  - Focusing on Faces
  - Interpretability of Face CNN
  - Action Units for Personality Traits Prediction
- Conclusions

# Outline

- Introduction
- Related Work
- Experiments
  - Images + audio vs Images for personality trait regression
  - Finding Discriminative Regions in video frames
  - Focusing on Faces
  - Interpretability of Face CNN
  - Action Units for Personality Traits Prediction
- Conclusions

#### Introduction

- Problem: Automatic apparent personality trait inference
  - Big Five apparent personality traits
- Approach: Interpret CNN models
  - What internal representations emerge?
  - What image regions are more discriminative?



#### Introduction

- Challenge: First Impressions dataset
  - Most recent and large database for apparent personality trait estimation
  - $\circ$  10,000 video clips
  - Video frames, audio and captions available
  - Big Five personality traits annotated in a continuous 0-1 scale



# Outline

- Introduction
- Related Work
- Experiments
  - Images + audio vs Images for personality trait regression
  - Finding Discriminative Regions in video frames
  - Focusing on Faces
  - Interpretability of Face CNN
  - Action Units for Personality Traits Prediction
- Conclusions

#### **Related Work**

- CNN models interpretability
  - Class Activation Map (CAM) [Zhou et al, CVPR'16]
    - Visualize class-specific discriminative regions



[Zhou et al, CVPR'16] "Learning deep features for discriminative localization."

#### **Related Work**

- Deep learning architectures for personality trait regression
  - Fully Convolutional Neural Network (Zhang et al, ECCVW'16)
    - Winner last edition on First Impressions challenge
    - This architecture has been used as reference
  - LSTM Recurrent Neural Network (Subramaniam et al, ECCVW'16)
  - Deep Residual Network (Güçlütürk et al, ECCVW'16)

[Zhang et al, ECCVW'16] "Deep bimodal regression for apparent personality analysis."

[Subramaniam et al, ECCVW'16] "Bi-model first impressions recognition using temporally ordered deep audio and stochastic visual features." [Güçlütürk et al, ECCVW'16] "Deep impression: audiovisual deep residual networks for multimodal apparent personality trait recognition."

#### **Related Work**

- Fully Convolutional Neural Network (Zhang et al, ECCVW'16)
  - 2 models (images and audio) + late fusion
  - Model for images: DAN+
    - Extension of DAN (Descriptor Aggregation Networks)
    - Pre-trained VGG-face model
    - Average and max pooling at 2 different layers
  - Model for audio
    - Regression model over log filter bank features



[Zhang et al, ECCVW'16] "Deep bimodal regression for apparent personality analysis."

# Outline

- Introduction
- Related Work

- Images + audio vs Images for personality trait regression
- Finding Discriminative Regions in video frames
- Focusing on Faces
- Interpretability of Face CNN
- Action Units for Personality Traits Prediction
- Conclusions

- 1. Images + audio vs Images for personality trait regression
  - Objective: Focusing only on image model interpretation
  - Accuracy of the models
    - Images (100 frames per video) + audio: 0.913
    - Only images (10 frames per video): 0.909

	Mean accuracy	Openness	Conscientiousness	Extraversion	Agreeableness	Neuroticism
img+audio	91.3	91.2	91.7	91.3	91.3	91.0
img	90.9	90.9	91.1	90.9	91.0	90.5

- 2. Finding Discriminative Regions in video frames
  - CAM (Class Activation Maps) is applied to the image model



Discriminative localization for 20 images with highest predicted value for agreeableness

**Class Activation Map** 

- 2. Finding Discriminative Regions in video frames
  - CAM (Class Activation Maps) is applied to the image model
  - Discriminative regions mainly on faces regions
  - Quantitative evaluation
    - Face detection algorithm
    - Overlap of face bbox and CAM regions

$$overlap = \frac{M_{face} \cap M_{CAM}}{M_{CAM}}$$



• Result: 72.80% of CAM regions have at least an overlap of 0.9 with the detected face

- 3. Focusing on Faces
  - Idea: Training the same architecture on cropped faces
  - Pre-processing:
    - Face region cropping
    - Eyes estimated localization for alignment
    - Image resize
  - Results:

	Mean accuracy	Openness	Conscientiousness	Extraversion	Agreeableness	Neuroticism
img	90.9	90.9	91.1	90.9	91.0	90.5
face	91.2	91.0	91.4	91.5	91.2	90.7

- 3. Focusing on Faces: Finding Discriminative Regions
  - CAM (Class Activation Maps) is applied to the image model



Discriminative localization for 20 images with highest predicted value for agreeableness

**Class Activation Map** 

- 4. Interpretability of Face CNN
  - Goal: Visualize whether semantic detectors emerge from the network
  - Methodology (based on Zhou et al, ICLR'15)
    - Visualization of images that produce the highest activation given a unit of a layer
    - Images are segmented using an estimated receptive field



[Zhou et al, ICLR'15] "Object detectors emerge in deep scene CNNs."

- 4. Interpretability of Face CNN
  - Result: Semantic regions such as eyes, nose and mouth emerge
  - Previous methodology: manual inspection
  - New approach: automatic identification of emerging semantic detectors
    - Images are aligned
    - Semantic regions are defined
    - Spatial histograms from highest activations localization are computed for each unit of the CNN architecture
    - The addition of the spatial histogram values for a specific semantic region is applied to identify semantic detectors



- 4. Interpretability of Face CNN
  - Eyebrow detectors



- 4. Interpretability of Face CNN
  - Eye detectors



- 4. Interpretability of Face CNN
  - Nose detectors



- 4. Interpretability of Face CNN
  - Mouth detectors



- 5. Action Units in Personality Traits Regression
  - Influence of shown emotion for personality trait
  - 17 Action Units (AU) from Facial Action Coding Systems
  - AU as 17-dimensional feature vector
  - Linear regressor trained on these feature vectors
  - Mean Accuracy: 88.6

	Mean accuracy
img	90.9
face	91.2
AUs	88.6

	Upper Face Action Units						
	AU 1	AU 2	AU 4	AU 5	AU 6	AU 7	
	10	100	10	100	6	-	
กร	Inner Brow Raiser	Outer Brow Raiser	Brow Lowerer	Upper Lid Raiser	Cheek Raiser	Lid Tightener	
-	*AU 41	*AU 42	*AU 43	AU 44	AU 45	AU 46	
	0	00	0	96	0	0	
	Lid Droop	Slit	Eyes Closed	Squint	Blink	Wink	
ĺ	Lower Face Action Units						
	AU 9	AU 10	AU 11	AU 12	AU 13	AU 14	
	1	1	31	10		1	
	Nose	Upper Lip	Nasolabial	Lip Corner	Cheek	Dimpler	
	Wrinkler	Raiser	Deepener	Puller	Puffer		
	AU 15	AU 16	AU 17	AU 18	AU 20	AU 22	
	12		3(1)		1	O/	
	Lip Corner	Lower Lip	Chin	Lip	Lip	Lip	
	Depressor	Depressor	Raiser	Puckerer	Stretcher	Funneler	
	AU 23	AU 24	*AU 25	*AU 26	*AU 27	AU 28	
			1	E)			
	Lip	Lip	Lips	Jaw	Mouth	Lip	
	Tightener	Pressor	Part	Drop	Stretch	Suck	

- 5. Emergence of Action Unit Detectors in Personality Traits Regression
  - Do AU detectors emerge from internal units of CNN model?
    - N frames with highest predicted intensity value for a given AU: {F<sub>AU</sub> }
    - N frames with highest activation for a given internal unit: {F<sub>unit</sub> }
    - Internal unit with highest intersection I<sub>max</sub> between {F<sub>AU</sub>} and {F<sub>unit</sub>} is identified
    - Probability p to obtain I<sub>max</sub> by chance is computed

• 5. Emergence of Action Unit Detectors in Personality Traits Regression



# Outline

- Introduction
- Related Work
- Experiments
  - Images + audio vs Images for personality trait regression
  - Finding Discriminative Regions in video frames
  - Focusing on Faces
  - Interpretability of Face CNN
  - Action Units for Personality Traits Prediction
- Conclusions

### Conclusions

- Interpretability of deep learning models for apparent personality trait inference
- Facial information was found to play a key role from discriminative region visualization
- Facial part detectors automatically emerged from last layers with no supervision provided on this task
- Influence of emotional information on trait prediction with the use of Action Units was explored



- Action Units for Personality Traits Prediction
  - Influence of shown emotion for personality trait inference
  - 17 Action Units (AU) from Facial Action Coding Systems
  - Do AU detectors emerge from internal units of CNN model?
    - N frames with highest predicted intensity value for a given AU: {F<sub>AU</sub> }
    - N frames with highest activation for a given internal unit: {F<sub>unit</sub> }
    - Internal unit with highest intersection I<sub>max</sub> between {F<sub>AU</sub>} and {F<sub>unit</sub>} is identified
    - Probability p to obtain I<sub>max</sub> by chance is computed

$$p = P(intersection \ge k) = 1 - \left(1 - \sum_{i=k}^{N} \frac{\binom{F-N}{N-i}\binom{N}{i}}{\binom{F}{N}}\right)^{U}$$

- Interpretability of Face CNN
  - Spatial histograms of the most frequent activation locations for each convolutional layer

