# DisplaceNet: Recognising Displaced People from Images by Exploiting Dominance Level

Grigorios Kalliatakis   Shoaib Ehsan   Maria Fasli   Klaus McDonald-Maier
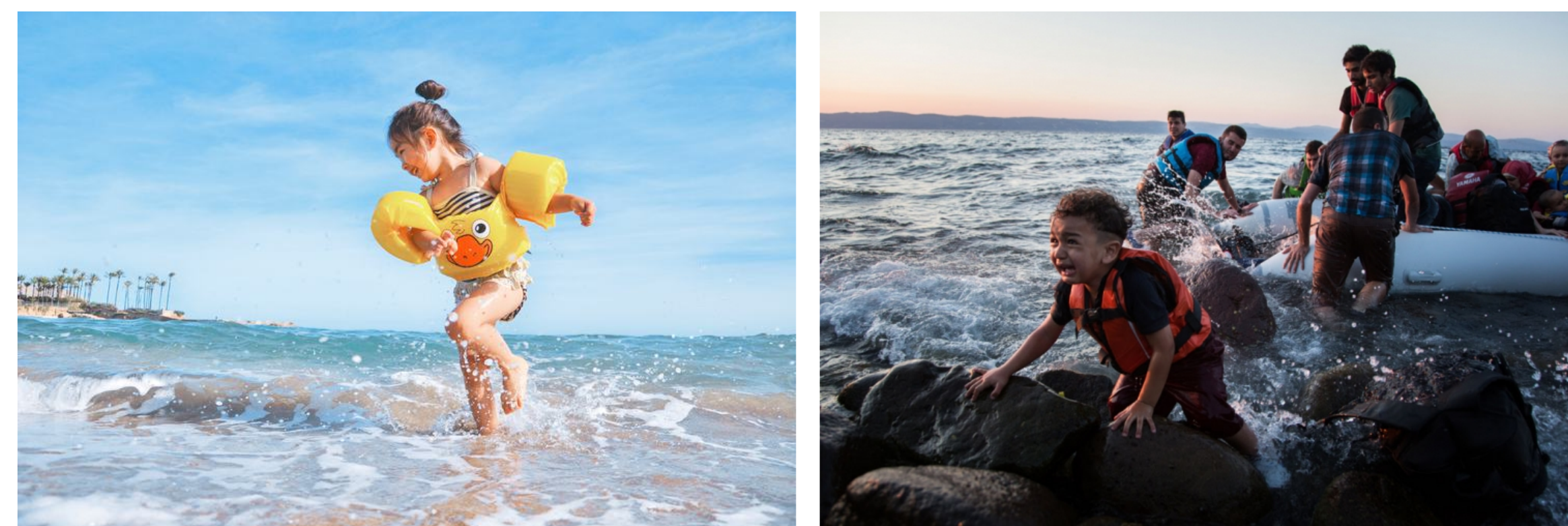
University of Essex, UK

University of Essex

The Human Rights, Big Data and Technology Project

## OVERVIEW

**Objective:** Inference of potential displaced people from real-world images.

**Motivations:**

- **68.5 million** forcibly displaced individuals worldwide - roughly equivalent to the entire UK population being forced to flee their homes.
- Traditional methods for human-rights-related image analysis require **manual labour** by human rights analysts and advocates.
- Computer vision can help **automate parts** of this process and turn recognition of displaced populations into a powerful and cost-effective application that could improve humanitarian responses.

## PROBLEM FORMULATION

**Can you label the images below as either displaced people or non-displaced people ?**
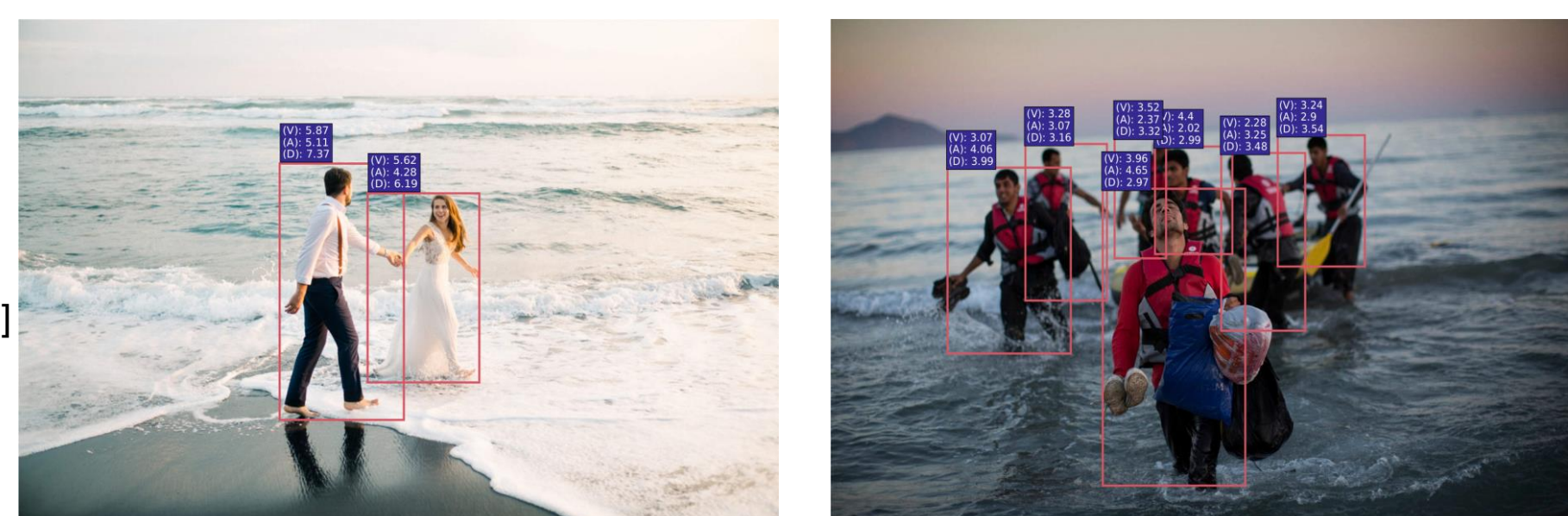Try to label them from the inference results of object detection and/or scene recognition.



**Main Idea:** A person's control level of a situation can be a notifying difference between the encoded visual content of an image that depicts a non-violent situation and the encoded visual content of an image displaying displaced people.
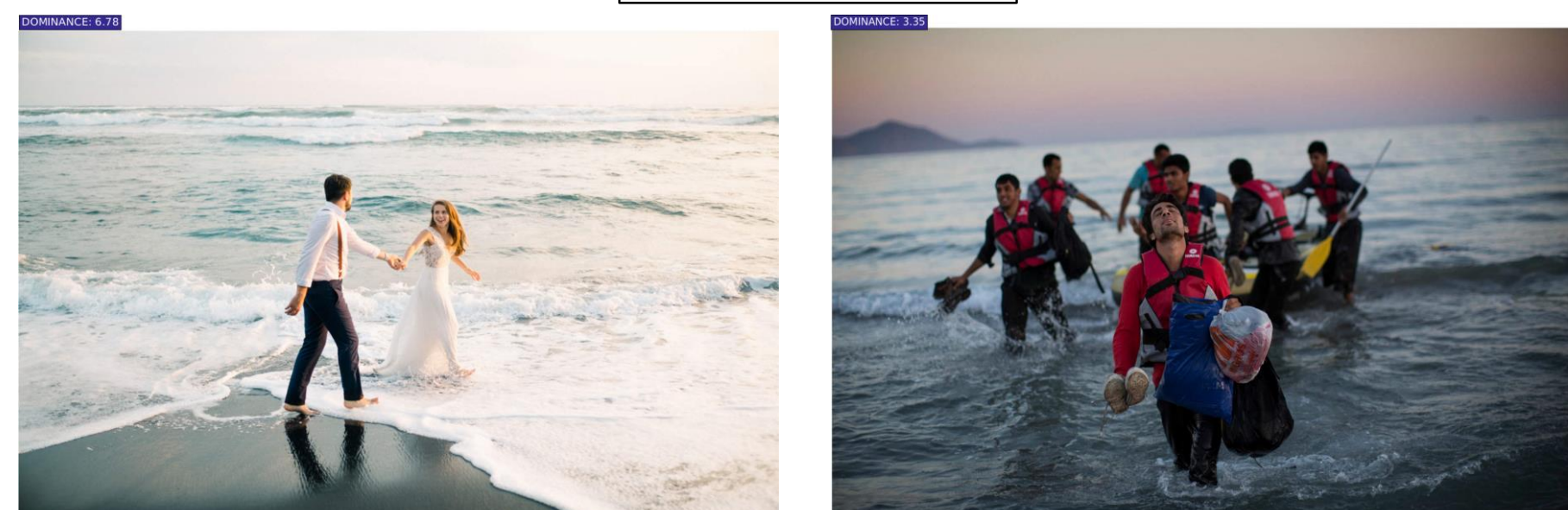
## OVERALL DOMINANCE SCORE

1. Combine the person bounding box with the information present in scene context similar to [1]–to recognise emotions expressed in continuous dimensions *Valence*, *Arousal* and *Dominance*
2. Introduce the **overall dominance score** that characterises an entire image based on all individuals' control level of the situation



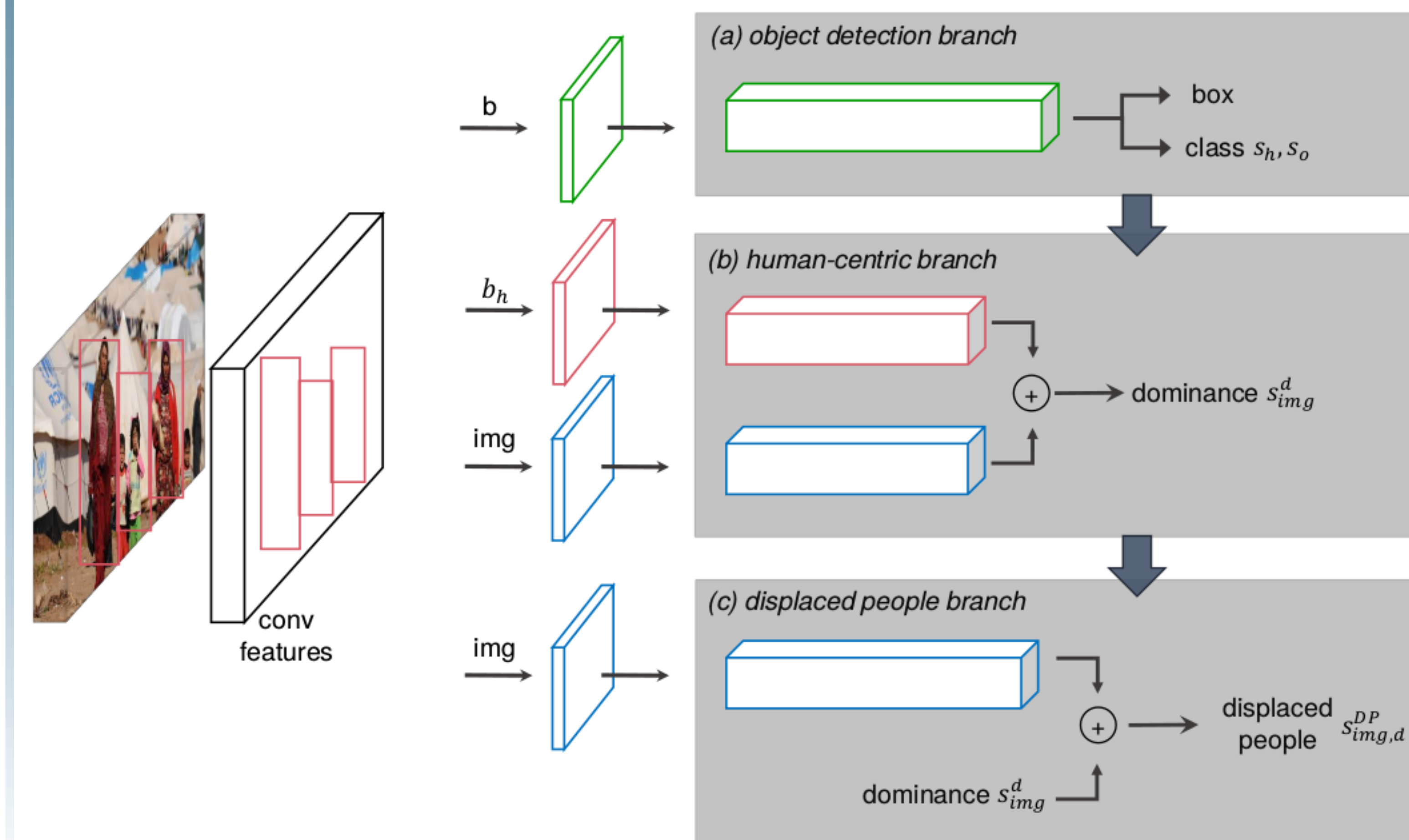Describe emotions using 3 numerical dimensions **Valence, Arousal & Dominance** [1]

$$s_{img}^d = \frac{1}{n} \sum_{i=1}^{n} s_{h,img}^d$$

Characterise the entire image using the proposed **Overall Dominance Score** $s_{img}^d$

## METHOD

**Model Architecture:**



**Proposed Solution:** Extend typical image classification by assigning a triplet score $s_{img,d}^{DP}$ to pairs of candidate human boxes $b_h$ and the displaced people category

$$s_{img,d}^{DP} = s_h \cdot s_{h,img}^d \cdot s_{img}^{DP}$$

**Components:**

**Object Detection Branch:** Localise the boxes containing a human $b_h$ and the object of interaction $b_o$ using RetinaNet [2].

**Human-centric Branch:**
- VAD score for each $b_h$.
- Dominance score $s_{img}^d$ that characterises **entire** image.
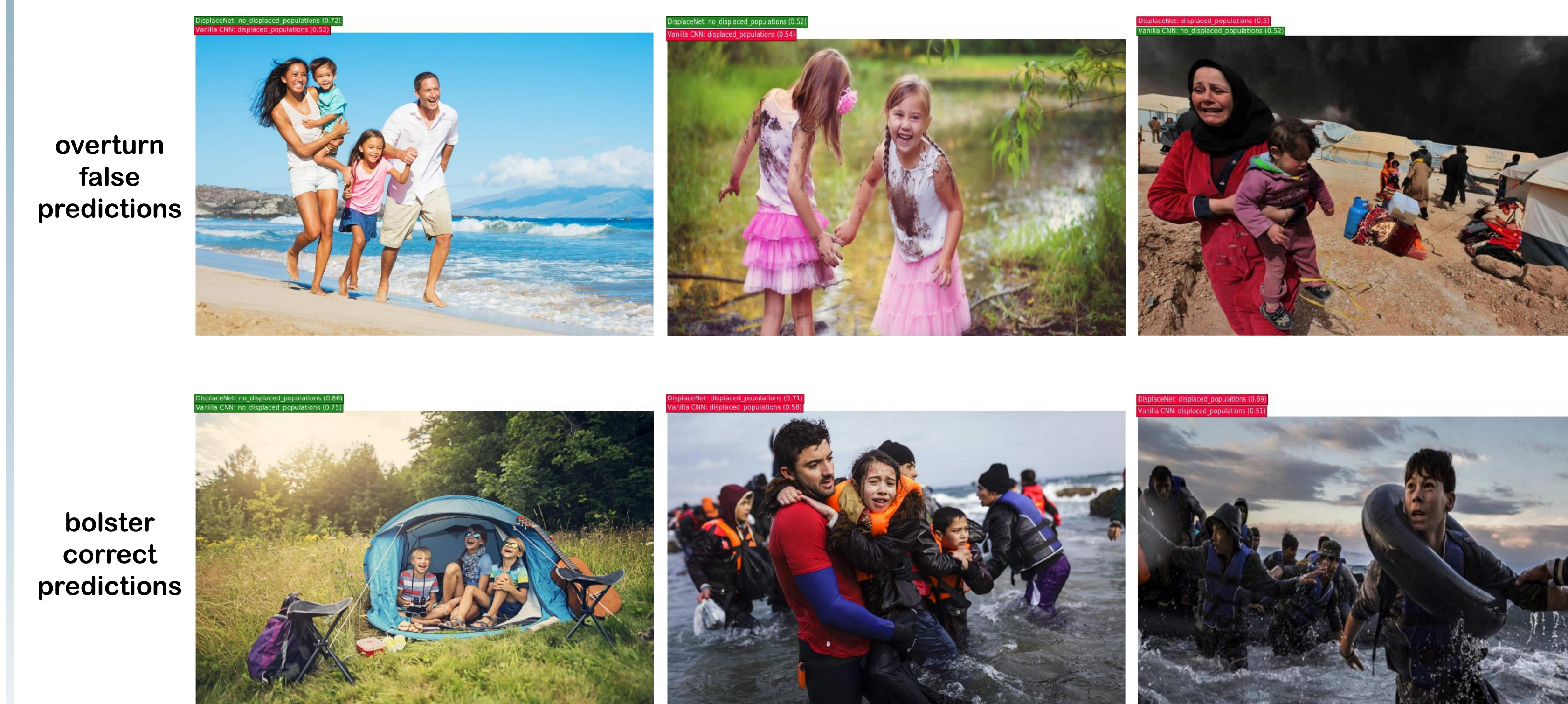
**Displaced People Branch:**
- Classification score for input image.
- Re-adjust classification score based on $s_{img}^d$.

## EXPERIMENTS & RESULTS

**Quantitative Comparison of vanilla CNNs* & DisplaceNet:**

| backbone network | layers fine-tuned | vanilla CNN | | DisplaceNet | |
|---|---|---|---|---|---|
| | | Top-1 acc. | Coverage | Top-1 acc. | **Coverage** |
| VGG16 | | 58% | 0% | 54% | **3%** |
| VGG19 | 1 | 69% | 3% | 60% | **6%** |
| ResNet50 | | 60% | 0% | 55% | **4%** |
| VGG16 | | 63% | 43% | 63% | **49%** |
| VGG19 | 2 | 77% | 54% | 74% | **58%** |
| ResNet50 | | 42% | 1% | 38% | **5%** |
| **mean** | - | 61.5% | 16.83% | 57.33% | **20.83%** |

**Qualitative Results:**



overturn false predictions

bolster correct predictions

*image classification using solely fine-tuning without any other modification

## DATASET & METRICS

- Two-class subset of Human Rights Archive Dataset [3].
- Use of *coverage*–proportion of a dataset for which a classifier is able to produce a prediction–as a realistic performance metric.
- DisplaceNet refuses to classify an input $x$, whenever the probability of the output sequence $p(y|x) < t$ for some confidence threshold $t = 0.85$.

## REFERENCES

[1] Emotion Recognition in Context [Kosti *et al.*, CVPR17]

[2] Focal loss for dense object detection [Lin *et al.*, ICCV17]

[3] Exploring Object-Centric and Scene-Centric CNN Features and Their Complementarity for Human Rights Violations Recognition in Images [Kalliatakis *et al.*, IEEE Access 2019]

## ACKNOWLEDGEMENTS