# DEEP NEURAL NETWORKS FOR OBJECT DETECTION:

## A Critical Analysis

Mohith Naga Adithya Vasamsetti

MXV00070@UCMO.EDU

University of Central Missouri, Lee's Summit

## I. INTRODUCTION:

Advancements in image object detection have significantly progressed over time, with deep neural networks (DNNs) at the forefront of these developments. This article provides a thorough examination of the research conducted by Christian Szegedy, Alexander Toshev, and Dumitru Erhan at Google Inc., entitled "Deep Neural Networks for Object Detection," introducing a fresh method for object detection utilizing DNNs. The goal of the paper is to tackle both the categorization and the accurate positioning of objects in images, suggesting a straightforward yet effective approach that establishes new standards in the field.

## II. SHORT SUMMARY:

The paper builds on the success of DNNs in image classification and extends their application to object detection. The authors propose treating object detection as a regression problem to object bounding box masks, utilizing a multi-scale inference procedure to produce high-resolution object detections efficiently. Their approach involves using a DNN-based regression to output binary masks for object bounding boxes and applying a simple bounding box inference for detection. They demonstrate state-of-the-art performance on the Pascal VOC dataset, showcasing the effectiveness of their method in achieving precise object localization.

## III. CRITICAL ANALYSIS:

1. **Strengths**
   - **Innovative Approach**: The paper introduces a novel way of using DNNs for object detection by framing it as a regression problem. This method bypasses the need for complex, manually engineered features and allows the network to learn robust object representations autonomously.

- **Multi-Scale Inference**: The multi-scale inference procedure is a key strength of the paper. By applying the DNN-based regression across different scales and image crops, the method efficiently handles objects of varying sizes and achieves high-resolution detection with low computational cost.

- **State-of-the-Art Performance**: The authors validate their approach with impressive results on the Pascal VOC challenge, demonstrating that their method outperforms traditional object detection systems like the Deformable Part-based Model (DPM).

- **Simplicity and Generalization**: The proposed method is straightforward to implement and does not require hand-designed models for capturing parts and their relations. This simplicity allows for easy applicability to a wide range of object classes, from rigid to deformable objects.

2. **Weaknesses**

- **Computational Resources**: While the paper emphasizes the efficiency of their method, it still relies on substantial computational resources, especially for training the DNNs. This requirement might limit the method's applicability in resource-constrained environments.

- **Fixed Output Dimension**: The network's fixed output dimension for binary masks could pose limitations in accurately detecting very small or very large objects within high-resolution images. Although the multi-scale approach mitigates this to some extent, it might not completely resolve the issue.

- **Complex Training Procedure**: The training procedure involves generating thousands of samples from each image, which can be computationally expensive and time-consuming. Additionally, balancing the positive and negative samples effectively is crucial for optimal performance.

3. **Potential Improvements**

- **Adaptive Output Resolution**: Future work could explore adaptive output resolutions based on the input image size and object dimensions, potentially enhancing detection accuracy for objects of varying scales.

- **Efficiency Enhancements**: Optimizing the training procedure and exploring methods to reduce computational overhead without compromising accuracy could broaden the method's applicability in practical scenarios.

- **Integration with Real-Time Systems**: Enhancing the method for real-time object detection applications, such as in autonomous vehicles or surveillance systems, could be a valuable direction for future research.

## IV. SEQUENCE OF DETECTION & WORKING:

**DNN-based Detection:**

- The core of our approach is a DNN-based regression towards an object mask. This model can generate masks for the full object as well as portions of the object. A single DNN regression can produce masks of multiple objects in an image. We apply the DNN localizer on a small set of large sub-windows to further increase localization precision.

**Detection as DNN Regression:**

- Our network is based on a convolutional DNN architecture with seven layers, the first five being convolutional and the last two fully connected. We replace the final softmax classifier with a regression layer that generates an object binary mask. The network is trained by minimizing the L2 error for predicting a ground truth mask.

**Precise Object Localization via DNN-generated Masks:**

- The approach generates high-quality masks, but challenges remain, such as disambiguating objects placed next to each other and dealing with the limits in output resolution. To address these issues, we generate multiple masks for robust localization and apply a multi-scale refinement of the DNN localizer.

**Multiple Masks for Robust Localization:**

- They generate several masks for each object, predicting the full object box and four halves (bottom, top, left, and right). These over-complete predictions help reduce uncertainty and disambiguate multiple objects placed next to each other.
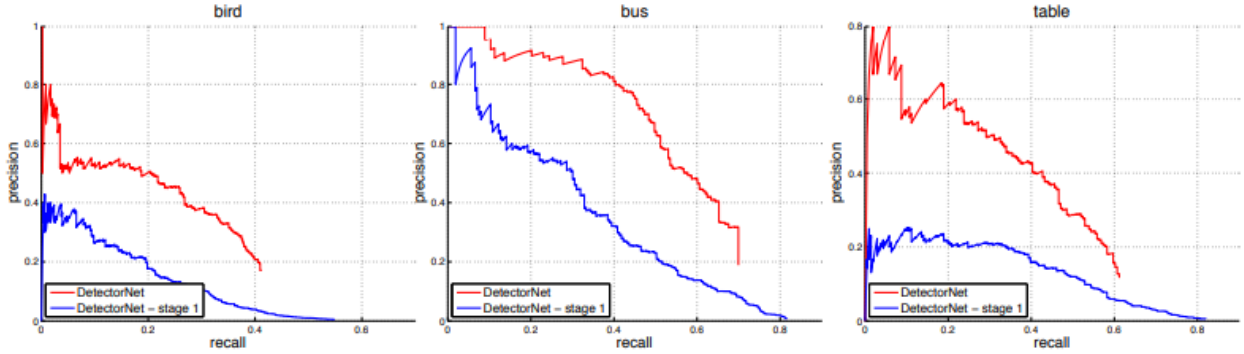
**Object Localization from DNN Output:**

- They estimate bounding boxes for each image by rescaling the binary masks to the original image size and inferring boxes with the highest scores. We use an agreement score to measure the consistency of each bounding box with the masks, performing an exhaustive search to find the best matches. Non-maximum suppression and a DNN classifier further refine the detections.

**Multi-scale Refinement of DNN Localizer:**

- They address output resolution limitations by applying the DNN localizer at several scales and large sub-windows, followed by a refinement step. This process involves merging masks generated at different scales and applying the DNN localizer on the top inferred bounding boxes, increasing detection precision.

## V. RESULTS:



Precision recall curves of DetectorNet after the first stage and after the refinement.

## V. CONCLUSION:

The research paper authored by Szegedy, Toshev, and Erhan signifies a notable progress in the realm of object detection through the utilization of deep neural networks. The authors have created a method that achieves top performance by treating object detection as a regression problem and using a multi-scale inference method, demonstrating simplicity and adaptability. Although facing some computational hurdles, the suggested method creates opportunities for precise and effective object detection in various contexts. Further enhancement of the method's practicality in real-world situations could be achieved through future research on adaptive output resolutions and computational optimizations.