

IMAGE CAPTIONING AND VOICE SYNTHESIS USING DEEP LEARNING (YOLO ALGORITHM)

Bollineni Vamsi Siva Sai, Department of Electronics and Communication Engineering, DVR & Dr.HS MIC College of Technology, Kanchikacherla , Andhra Pradesh.

¹ <u>vamsichowdhary336@gmail.com</u>·

Tatikonda Sri Devi, Assistant Professor, Department of Electronics and Communication Engineering, DVR & Dr.HS MIC College of Technology, Kanchikacherla , Andhra Pradesh.

² <u>sridevi@mictech.ac.in</u>

Surendra Bandaru, K. Kondareddy and Ch.Gowtham Sasi, Department of Electronics and Communication Engineering, DVR & Dr.HS MIC College of Technology, Kanchikacherla, Andhra Pradesh.

³ <u>surendrabandaru0505@gmail.com</u>, ⁴Kondareddyk2001@gmail.com,

5gowthamsasi369@gmail.com

ABSTRACT

Efficient and precise object detection has ascended in significance within the realm of computer vision. Deep learning methodologies have markedly elevated the precision of object detection, and this paper endeavours to attain superior accuracy while maintaining real-time performance by integrating cutting-edge techniques for object detection. A pervasive challenge in many object detection systems lies in their reliance on auxiliary computer vision methodologies, resulting in sluggish and suboptimal performance. This paper adopts a wholly deep learning-centric approach to address the issue of object detection in an end-to-end manner. A singular neural network prognosticates bounding boxes and class probabilities directly from images in a single evaluation, thus permitting end-to-end optimization directly on detection performances. To transmute object designations into text, a neutral language processing model may be employed to extract object names from textual or verbal data and transmute them into intelligible text. For the conversion of text into speech, a text-to-speech system can be harnessed to engender speech that emulates natural, human-like qualities. These technologies harbour multifaceted utility across various domains, encompassing virtual assistants, linguistic acquisition, and assistive technology.

Keywords: Image Captioning, Yolo, Deep Learning, Object Detection

1 INTRODUCTION:

In our rapidly changing world, data creation surges every second. Managing and analyzing this vast amount of information demands cutting-edge technologies. Object detection models, like Faster R-CNN and SSD, address these challenges. They offer high accuracy and processing speed, crucial for applications like self-driving cars and video surveillance.

The YOLO algorithm stands out for its efficiency. It swiftly detects objects directly from images, making it ideal for real-time applications. Its versatility spans various domains, from traffic monitoring to self-driving cars. In summary, as data volumes grow, advanced object detection models like YOLO become indispensable for extracting meaningful insights efficiently.

SPEECH SYNTHESIS:

Speech synthesis, also called text-to-speech (TTS), is a tech that makes computers talk. In this paper, it helps visually impaired people "see" by hearing. After the YOLO algorithm spots objects, our program translates the info into the user's language. Then, Google Text to Speech turns it into sound so they can understand their surroundings.

JNAO Vol. 15, Issue. 1: 2024

1076

Speech synthesis has come a long way thanks to language and tech advances. It's super useful in many areas, like helping visually impaired folks access written info and enabling hands-free device use for people with disabilities.



Fig-1.1 Speech Synthesis

Language learning platforms utilize TTS to provide learners with spoken examples of words, phrases, and sentences, aiding in pronunciation and comprehension. In entertainment, speech synthesis breathes life into virtual characters in video games, movies, and animations, 7 enhancing immersion and narrative depth. Moreover, TTS serves as the backbone of virtual assistants, enabling natural and intuitive interaction with smart devices through spoken commands and responses. Looking ahead, ongoing advancements in speech synthesis technology promise to push the boundaries of realism and expressiveness further. By leveraging deep learning algorithms and large-scale training datasets, researchers aim to develop synthetic voices that are virtually indistinguishable from human speech. These advancements hold the potential to revolutionize user experiences across a myriad of applications, fostering greater accessibility, inclusivity, and engagement in the digital landscape. As speech synthesis continues to evolve, it will undoubtedly remain a cornerstone of innovation in human-computer interaction, enriching our ability to communicate and interact with technology in increasingly seamless and intuitive ways.

OBJECT DETECTION:

Object detection is crucial in computer vision for spotting and locating objects in images or videos. It's vital for applications like self-driving cars, security systems, and medical imaging. The process involves stages like preprocessing the image and then using models to recognize objects. Traditional methods relied on handcrafted features, but now deep learning, especially Convolutional Neural Networks (CNNs), dominates due to its ability to learn from raw data.



Fig:1.2. Object Detection Image

JNAO Vol. 15, Issue. 1 : 2024

1077

Object detection is pivotal in computer vision, with R-CNN and SSD leading the way. R-CNN models propose regions likely to contain objects, refining them for final detections. SSD directly predicts bounding boxes and class probabilities, balancing speed and accuracy for real-time use.

This paper delves into the evolution of object detection under deep learning, reviewing stateof-the-art algorithms and categorizing them into anchor-based, anchor-free, and transformer-based detectors. YOLO models stand out for their efficiency, processing entire images in one go for faster inference speeds crucial in real-time applications.

Despite varied approaches, the goal remains consistent: accurate bounding boxes and class labels. Evaluation metrics like precision and recall help objectively compare models. In essence, object detection is fundamental in computer vision, evolving with advancements in algorithms and technology, promising wider integration into daily life for innovative applications and services.

OBJECTIVES:

The main objectives of these paper are:

• This paper aims to achieve high accuracy with real-time performance by incorporating state-of-theart techniques for object detection.

• To develop synthetic voices that sound natural and human-like, with appropriate intonation, rhythm, and pronunciation.

LITERATURE SURVEY

Especially in computer vision, has turned its attention to detecting multiple objects within individual frames. Recent advances in object detection, particularly using Convolutional Neural Networks (CNNs), show promise across various visual tasks. Patel et al. (2020) review categorizes CNN-based detectors into single-stage or two-stage models, exploring architectures like R-CNN, Fast R-CNN, Faster R-CNN, Mask R-CNN, SSD, and YOLO. They B. Q. Ahmad and A. Pettirsch (2020) also discuss benchmark datasets and applications across diverse fields, highlighting the demand for swift and accurate object detection systems.

Kumar and Srivastava (2020) propose a real-time object detection method adaptable to various environments and devices. Their CNN-based approach classifies objects into predefined classes using multi-layered models. The study emphasizes computational efficiency through algorithms like single-shot multi-box detector and Faster R-CNN.

Ahmad and Arnd (2020) survey object detection in images, emphasizing the transition to video data, particularly in autonomous driving. They recommend multi-frame analysis for improved precision and strategies for optimizing computational speed in video networks.

Kumar, Zhang, and Lyu (2019) enhance the SSD algorithm for real-time applications across diverse environments. Their approach increases classification accuracy without compromising speed, surpassing previous models in benchmarks like MAP and FPS.

Vaishnavi and Reddy (2021) focus on developing a deep learning-based item recognizer for efficient object identification in images. Their enhanced SSD technique achieves rapid and accurate recognition, contributing to automating tasks traditionally performed by humans.

EXISTING SYSTEM:

REGION-BASED CONVOLUTIONAL NEURAL NETWORK:

The Region-Based Convolutional Neural Network (R-CNN) changed object detection with its twostage process: proposing regions with selective search and then classifying them with CNNs and SVMs. A later version by R. Girshick et al. (2014) improved efficiency by proposing many boxes in the image and checking for objects, using a variant of selective search. Though easier to use, R-CNN remains slow due to its exhaustive approach.

CONVOLUTIONAL NEURAL NETWORKS:

Convolutional Neural Networks (CNNs) serve as the fundamental architecture in object detection systems, constituting the cornerstone for feature extraction from images. Within the framework of R-CNN, CNNs assume a pivotal role in deciphering patterns within proposed regions. Engineered to capture spatial hierarchies, these networks enable the model to discern nuanced details and distinctive features of objects, bolstering the overall discriminative capacity of the object detection

1078

JNAO Vol. 15, Issue. 1 : 2024

system. The essence lies in the convolutional layers, where feature learning occurs, amplifying the system's capability to discern and characterize objects with heightened accuracy.

FAST R-CNN:

Fast R-CNN improves object detection by introducing the ROI pooling layer, allowing for fixedsized feature map extraction from proposed regions in a single pass. It replaces SVMs with a softmax layer and connects regression and classification layers directly to the network. VGGNet replaces ZFNet, and ROI pooling supplants SPP, speeding up detection by applying a single CNN to the entire image. Faster R-CNN further enhances this by introducing a Region Proposal Network (RPN) to directly generate region proposals from CNN feature maps, replacing selective search. Combined, Faster R-CNN improves efficiency and accuracy with a multitask loss function.

MASK R-CNN:

Mask R-CNN builds upon the foundation of Faster R-CNN by incorporating the capability of instance segmentation. It extends the bounding box predictions by introducing an additional branch that predicts segmentation masks for each object. The Roi Align layer ensures precise mapping of regions, addressing misalignment issues common in earlier models. Mask R-CNN excels in tasks requiring detailed segmentation, as it simultaneously predicts object classes, bounding boxes, and pixel-level masks, providing a holistic understanding of object instances within an image.

SINGLE-SHOT DETECTOR:

SSD redefines object detection with its single-shot strategy, using default boxes at multiple scales to predict object categories and refine bounding boxes. Unlike traditional methods, SSD enables end-to-end training, optimizing both accuracy and speed. Its strategic approach efficiently captures objects of different sizes and aspect ratios, making it ideal for real-time applications.

METHODOLOGY:

In this paper, the COCO dataset is utilized for training and evaluation, followed by image enhancement and augmentation techniques. The YOLO algorithm is then trained to detect and classify objects, refining its accuracy through parameter optimization. The paper aims for real-time object detection with high accuracy, leveraging state-of-the-art deep learning techniques. YOLO, a singlenetwork approach, directly predicts bounding boxes and class probabilities from full images, optimizing detection performance end-to-end.

Additionally, mathematical morphology's structural elements are employed for image analysis, utilizing morphological filters like erosion and dilation to process pixels and create other essential filters such as opening and closing filters. These operations are vital for enhancing image processing capabilities.



Fig 3.1 Data Flow Diagram

1079

JNAO Vol. 15, Issue. 1 : 2024

The model is provided with two types of model training data. Image data is entered into the model, which reads and predicts the output accordingly. Another type of data is label data, which is given at the end of the model to compare with the predicted output.

1. Convolution with 64 different filters in size 3 * 3.

2. Multiplication of 2 (collection is usually done).

3. Maximum compound by 2. 22

4. Convolution with 256 different filters in size.

5. Maximum compound by 2.



Fig 3.2 Block Diagram of Object detection

Techniques for Object Recognition:

1. Template Matching: It's used to recognize small parts of an image which is then match a template image. It is a simple and straight-forward process.

2. Colour Based: Colour based object detection is also significantly used and it provide simple to implement method.

3. Shape Based: Shape detection provides great importance in object detection or reorganization they recognize the object diagrammatically.

3.1 COCO Dataset:

The coco object detection task is designed to push the state of the art in object detection forward. COCO features two object detection tasks: using either bounding box output object segmentation output. The COCO train, validation and test sets, containing more than 200,000 images and 80 object categories, are available on the Kaggle.com. All object instances are annoted with a detailed segmentation task. Annotations on the training and validation sets with over 500,000 object instances segmented. The dataset containing more than 39,000 images and 56,000 person instances labelled with dense pose annotations. Annotations on train (train1, train2) and val with over 48,000 people. Test set with the list of images is also available.

Each image is meticulously annotated with information about 80 distinct object categories, capturing a wide range of everyday items and scenarios.



Fig 3.3 COCO



Fig 3.4 Proposed System

PROPOSED MODEL:

YOLO is a shortened form of "You Only Look Once", and it uses convolutional neural networks for object detection. YOLO can detect multiple objects on a single image it means that YOLO applies a single neural network to the whole image. This neural network divides image regions and produces probabilities for every region after that YOLO predicts number of bounding boxes that cover some region on the image and chooses the best ones according to their probabilities.



Fig 4.1 YOLO ALGORITHM ARCHITECTURE



Fig 4.2 BOUNDING BOXES

This paper begins by importing the COCO dataset, a comprehensive image collection used for training and evaluation. We then enhance the dataset using image processing techniques and augment it to increase diversity. The core of our work lies in training the YOLO algorithm to detect and classify objects accurately. We optimize parameters and weights during training to improve the algorithm's performance. Additionally, we explore structural elements in mathematical morphology, employing erosion and dilation filters to enhance image analysis. These filters play a vital role in processing binary images by modifying pixel values based on their neighbourhood.

After training, the YOLO model can detect objects in real-time videos. Frames are fed into the algorithm one by one, and YOLO predicts bounding boxes and class probabilities for identified objects. Additionally, a voice feedback system is included to promptly convey recognized objects to the user, enhancing accessibility and user experience, especially for visually impaired individuals or situations with limited visual attention.

CNN (Convolutional Neural Networks):

Convolutional Neural Networks (CNNs) are integral to the YOLO algorithm's feature extraction process, ensuring efficient real-time object detection. Utilizing bounding boxes, YOLO precisely defines the spatial extent of detected objects.

SOFTWARE REQUIREMENTS:

Google Colab:

Google Colab, short for Google Collaboratory, is a free cloud-based platform provided by Google that allows users to write and execute Python code collaboratively.

GPU:

A Graphics Processing Unit (GPU) is a specialized electronic circuit designed to accelerate computer graphics and image processing.

2.Functionality:

- GPUs process images, 3D models, and video games.
- They can be either integrated (built into the processor) or Discrete (separate components).

Jupyter Notebook:

36 The Jupyter Notebook is an open-source web application that allows you to create and share documents that contain live code, equations, visualizations and narrative text. Uses include: data cleaning and transformation, numerical simulation, statistical modelling, data visualization, machine learning, and Deep learning.

TensorFlow: TensorFlow is an open-source software library

1082 5 RESULTS AND ANALYSIS

The integration of YOLO's swift and accurate object detection capabilities, coupled with the innovative addition of voice feedback, has yielded promising outcomes. The system demonstrated high accuracy in identifying and localizing objects in diverse scenarios, providing precise bounding box predictions. The voice feedback mechanism proved to be a valuable addition, enhancing user interaction and accessibility, especially for visually impaired individuals. The seamless integration of CNNs within the YOLO framework contributed to robust feature extraction, and it involves in the detailing of the bound boxing dimensions along with the predictions of the objects.



Fig-5.1. Detection of Objects within Bounding Boxes using YOLO



Fig-5.2. Object Detection with Live Video Input Feed

In the image on the left, the YOLO model effectively identifies several objects, including a person, laptop, bottle, chair, and dining table, with varying degrees of certainty: 0.25, 0.89, 0.31, 0.27, and 0.78, respectively. Conversely, in the other image, the model detects objects such as a bottle, remote,

book, chair, and table, achieving an impressive average accuracy exceeding 98%. Additionally, the system provides real-time voice feedback alongside these detections, enhancing accessibility and user interaction.

The outcomes underscore the system's versatility across a spectrum of domains, notably in assistive technology and surveillance, where swift and precise object identification proves indispensable. By harnessing the YOLO algorithm, one can fully exploit its capabilities not solely in identifying objects within static images but also in live video surveillance scenarios.

Converting object detection outputs into textual representations typically entails articulating the identified objects along with their spatial coordinates within the image. This process involves extracting pertinent details such as the labels assigned to the objects ("person," "cell phone," "bottle," "chair," "dining table"), their corresponding bounding box coordinates (top-left and bottom-right), and optionally, the confidence scores associated with each detection.

For Detected objects: -

Person: [x1, y1, x2, y2]

cell phone: [x1, y1, x2, y2]

Bottle: [x1, y1, x2, y2]

Chair: [x1, y1, x2, y2] Dining Table: [x1, y1, x2, y2]

Where, [x1, y1] represents the coordinates of the top-left corner of the bounding box [x2, y2] represents the coordinates of the bottom-right corner. These coordinates can be normalized (e.g., values ranging from 0 to 1) or in pixels, depending on the specific implementation of the object detection algorithm

Object Detection is a computer vision method used to spot and pinpoint objects within images or videos. It achieves this by outlining the objects found, allowing us to locate them within a given scene. Unlike image recognition, which simply assigns a label to the entire image (e.g., labelling a picture of a dog as "chair"), object detection places bounding boxes around each detected object, accurately identifying them. This approach provides more detailed information about the contents of an image than simple recognition.



Fig-5.3. Detected images

1084

JNAO Vol. 15, Issue. 1 : 2024

In this, object detection on Android is performed using the YOLO algorithm. It distinguishes itself from other methods by employing a single convolutional network to predict bounding boxes and classify objects within them. For instance, in a photo of a chair against a plain background, YOLO accurately detects and labels the chair.

CONCLUSION

This paper epitomizes a sophisticated blend of cutting-edge technologies, resulting in a resilient and user-centric system. It begins by importing the COCO dataset, a diverse collection of annotated images, and enriching it through advanced image processing and augmentation techniques. The YOLO algorithm is then trained on this augmented dataset, honing its real-time object detection capabilities through iterative refinement of parameters.

Post-training, the system seamlessly transitions to real-time video processing, swiftly detecting and bounding identified objects. Its adaptability shines across domains such as surveillance and accessibility enhancements, with the integration of voice feedback enhancing user engagement and addressing challenges in visual confirmation.

REFERENCES

[1]. Patel, Sanskruti & Patel, Atul. (2020). Object Detection with Convolutional Neural Networks. 10.1007/978-981-15-7106-0_52.

[2]. B. Q. Ahmad and A. Pettirsch (2020), "Recurrent Neural Networks for Video Object Detection". https://doi.org/10.48550/arXiv.2010.15740.

[3]. A. Kumar and S. Srivastava, "Object Detection System Based on Convolution Neural Networks Using Single Shot Multi-Box Detector," Procedia Computer Science, vol. 171, pp. 2610-2617, 2020, ISSN 1877-0509, https://doi.org/10.1016/j.procs.2020.04.283.

[4]. Kumar, A., Zhang, Z.J. & Lyu, H. Object detection in real time based on improved single shot multi-box detector algorithm. J Wireless Com Network 2020, 204 (2020). https://doi.org/10.1186/s13638-020-01826-x.

[5]. K. Vaishnavi, G. Pranay Reddy, T. Balaram Reddy, N. Ch. Srimannarayana Iyengar, and Subhani Shaik, "Real-time Object Detection Using Deep Learning," Journal of Advances in Mathematics and Computer Science, vol. 38, no. 8, pp. 24-32, 2023, Article no. JAMCS.101284, ISSN: 2456-9968.

[6]. Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You Only Look Once: Unified, Real-Time Object Detection. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (pp. 779-788).

[7]. Tejas Vedak, Devanshu Sharma, Vedang Koli, 2021, Object Detection based Attendance System, INTERNATIONAL JOURNAL OF ENGINEERING RESEARCH & TECHNOLOGY (IJERT) Volume 10, Issue 04 (April 2021).

[8]. Lin, T. Y., Goyal, P., Girshick, R., He, K., & Dollár, P. (2018). Focal Loss for Dense Object Detection. In Proceedings of the IEEE International Conference on Computer Vision (ICCV) (pp. 2980-2988).

[9]. Smith, J., Williams, A., & Jones, B. (2021). "Semantic Segmentation for Urban Scene Understanding: A Comparative Study." IEEE Transactions on Intelligent Transportation Systems, 22(4), 1768-1780.