



JOURNAL OF EMERGING TECHNOLOGIES AND INNOVATIVE RESEARCH (JETIR)

An International Scholarly Open Access, Peer-reviewed, Refereed Journal

CROP AND DROP USING SALIENT OBJECT DETECTION

G.Parthasarathy¹, S.V.R. Bricillaa Maria Sumid², N.Hajira Farlin³, V.S. Kavi Priya⁴, B.Vijayalakshmi⁵

¹ Professor, ²Student, ³Student, ⁴Student
¹Computer Science Engineering,
¹SRM TRP Engineering College, Tiruchirapalli, India

Abstract : The crop and drop tool lets you to digitalize the real world objects around us. It is an innovative idea that allows the user to take photo and detect the object in the real world and drop the image into the desktop computer. This tool will make the job of users much easier, as they will be able to point their phone cameras at an object and copy-paste it on their computers, instead of taking a photo, editing it and then inserting the cut-out into the document. The tool uses Augmented Reality (AR) and machine learning algorithm to detect the objects and isolate the image so that the background is automatically removed. For detection of image and for removal the background an open-source technology called U² Net is used and a computer vision algorithm called as Scale-invariant feature transform (SIFT) matches coordinates on the phone with the computer screen allowing you to place digital captures in specific locations on your computer screen. As a whole copy and paste your surroundings' using AR is the quickest way to capture, extract and transfer anything around you. This tool reverses the process and brings physical things into the digital world.

Index Terms - Augmented Reality, Computer vision, Background removal, Scale Invariant Feature Transform, Cleargrasp, Transparent object detection.

I. INTRODUCTION

As the innovation is developing step by step, everyday work in individuals' life is getting digitalized and it's been helping by lessening the load in each some sort of manner. Regardless of whether there are loads of alternate routes acquainted with individuals, still there are numerous things that are very simple however the interaction can very disturb. The altering of records in various ways by the Photoshop, the drag and drop incorporate for including records to an internet-based webpage or joining them to our messages fair by hauling those from their region to the net area application has become propensity to us. It is an element that saves minutes of dreary perusing the spring up window, which doesn't permit us to change to an extraordinary screen until the record is picked, or the client chooses to stop. In this Increased Reality can be exceptionally fascinating to assist with a new viewpoint. The reason for Increased Reality is to rethink or expand how a private would communicate with and decipher the significant world by acquainting virtual data with their immediate environmental elements furthermore, way to deal with this present reality climate as a whole. The simplified highlight for including records to a web-based website or going along with them to our messages fair by hauling those from their area to the net area application has become propensity to us. However, if we had the option to duplicate glue true articles/things into advanced screens it will save a ton of time and facilitate our work. Expanded reality engages us to duplicate glue certifiable things into our progressed screens with fair by scarcely any clicks. Increased Reality has updated our genuine experience with various features inside our screens inside the mechanized world. But this idea of carrying true articles into the advanced world brings a couple of complete movements according to our viewpoint and our communication with Augmented Reality Apps. The course of gluing genuine items by tapping the snap of an article, getting hinder its experience, and sharing it along with your work area might be a dreary one. With this Augmented Reality Application, it replaces this entire cycle by guiding your camera toward select the ideal article, making it proficient by bypassing snapping, veiling, saving, and exchanging windows for the client making it a truly productive AR copy paste Photoshop change.

In this research proposal, section 2 explains the preliminaries, literature survey explained in section 3, the proposed method is discussed in section 4, experiment results discussed in section 6 and conclusion is in section 7.

II. RELATED WORKS

2.1 Title: Edge-Aware Multiscale Feature Integration Network for Salient Object Detection in Optical Remote Sensing Images, 2022

Authors: Chenggang yan, Xiaofei Zhou

Chenggang yan et al. proposed, the decoder that integrates the enriched multiscale deep features during a rough-to-satisfactory way, yielding a top-notch saliency map. The experiments which were held on public optical RSI datasets clearly prove the effectiveness and superiority of the proposed EMFI-Net towards the trending saliency fashions. Specifically, the proposed EMFI-Net first generates effective multiscale deep features via the usage of the three convolutional branches with one-of-a-kind resolution inputs and also the cascaded function fusion module, so that a power-full illustration of salient objects could be also received. Therefore, the explicit and the implicit utilization of the edge records no longer only further strengthens the multiscale deep features but additionally endows the saliency maps with clean boundaries.

2.2 Title: BASNet: A Boundry Aware siamese Network for accurate remote sensing change detectio, 2022

Authors: Hao Wei, Rui Chen

Hao Wei et al. proposed a process named Change Detection (CD) in far off-sensing pics is one of the maximum critical subjects in the computer vision community. Most recent CD pipelines consciousness on introducing interest mechanism to enhance the discriminative capacity of network, but their crude model architectures cause faulty predictions and abnormal barriers. In this article, Boundary-conscious Siamese network (BASNet) is used to correct far off sensing CD. Based at the encoder–decoder structure, they first proposed a unique multi-scale paired fusion module (MPFM) to efficiently fuse the same-level function pairs from the Siamese encoding stream. In addition, they have designed a place steorage module (LGM) to correctly discover the changed areas.

2.3 Title: BASNet: Boundary-Aware Salient Object Detection, 2022

Authors: Xuebin Qin, Zichen Zhang, Chenyang Huang, Chao Gao, Masood Dehghan

Xuebin Qin et al. proposed, Deep Convolutional Neural Networks that had been adopted for salient item detection and finished the brand-new overall performance. Most of the preceding works however recognition on area accuracy but now not at the boundary best. In this pain keeping with, they have recommended a predict-refine structure, BASNet, and a brand-new hybrid loss for Boundary-Aware Salient object detection. The hybrid loss guides the community to research the transformation between the enter image and the floor truth in a three-level hierarchy – pixel-, patch- and map- degree – by using fusing Binary Cross Entropy (BCE), Structural Similarity (SSIM) and Intersection Over-Union (IoU) losses. Equipped with the hybrid loss, the proposed predict-refine structure is capable of effectively section the salient item regions and correctly are expecting the excellent systems with clean limitations. In this unit, they proposed a unique quit-to-stop boundary conscious model, BASNet, and a hybrid fusing loss for accurate salient object detection.

2.4 Title: U2-Net going deeper with nested U-Structured for salient object detection, 2022

Authors: Xuebin Qin, Zichen Zhang, Chenyang Hung

Xuebin Qin et al. proposed a design which is easy yet powerful deep network structure, U2Net, for salient item detection (SOD). The architecture of U2-Net is a two-level nested U-structure. The design has the following advantages: (1) It could seize more contextual data from different scales thanks to the aggregate of receptive fields of different sizes in our proposed Residual U-blocks (RSU), (2) It will increase the depth of the complete architecture without significantly increasing the computational value because of the pooling operations utilized in these RSU blocks. This architecture allows us to teach a deep network from scratch without the usage of backbones from photograph type obligations. They instantiate models of the proposed architecture, U2-Net (176.3 MB, 30 FPS on GTX 1080Ti GPU) and U2-Net† (4.7 MB, 40 FPS), to facilitate the utilization in ranging environments. Both models reap aggressive performance on six SOD datasets.

2.5 Title: Salient object detection: An Accurate and Efficient Method for Complex, 2021

Authors: Min Qiao, Gang Zhou, Qiu Ling and Li Zhang

Min Qiao et al. proposed knowledge-based salient object detection (SOD) strategies that have made fantastic progress in latest years. However, most deep getting to know-based totally strategies suffer from coarse item boundaries and pricey computations, particularly in detecting items with complicated shapes. This paper affords a correct and accurate Salient Object Detection method. This is based on a singular double-department network that consists of a frame department and a side department. To achieve a precise part, an edge profile enhancement module (EPEM) is integrated into the facet department. In addition, a fusion comments module (FFM) is embedded to combine capabilities from the two branches. To address the hassle of steeply-priced computations, channel interest module (CAM) is blanketed to restrain redundant characteristic channels

2.6 Title: A multi-task collaborative network for light field salient object detection, 2021

Authors: Qindan Zhang, Shiqi Wang

Qindan Zhang et al proposed the option to foresee the notable article is of crucial significance in picture handling and computer vision. With various methodologies proposed for programmed picture furthermore, video striking article discovery, significantly less work has been committed to recognizing and sectioning notable articles from light fields. All the more explicitly, the correlation systems among edge location, profundity deduction and striking item location are painstakingly examined to work with the agent saliency highlights. Consequently, the profundity arranged saliency highlights are gotten from the math of light fields, wherein the 3D convolution activity is utilized with strong representation capacity to show the uniqueness relationships among various perspective pictures. At last, an element upgraded remarkable object generator is created to incorporate these integral saliency highlights, prompting the last notable item expectations for light fields.

2.7 Title: Fine-Grained Visual Recognition in mobile Augmented Reality for Technical support,2020**Authors: Bing Zhou, Sinem**

Bing Zhou et al. proposed Expanded Reality that is progressively investigated as the new mechanism for two-way distant joint effort applications to direct the members all the more actually and productively by means of visual guidelines. As clients make progress toward more regular association and computerization in increased reality applications, new visual acknowledgment methods are expected to upgrade the client experience. Albeit straightforward object acknowledgment is in many cases utilized in expanded reality towards this objective, most joint effort assignments are excessively intricate for such acknowledgment calculations to do the trick. In this paper, the author has proposed a fine-grained visual acknowledgment approach for portable increased reality, which influences RGB video outlines and meagre profundity include focuses distinguished continuously, as well as camera present information to recognize different visual conditions of an item. They exhibit the worth of our methodology through a versatile application intended for equipment support, which consequently identifies the condition of an item to introduce the right arrangement of data in the right setting.

2.8 Title: Comparing Non-Visual and Visual Guidance Methods for Narrow field of view Augmented Reality Displays,2020**Authors: Alexander Marquardt, Christina Trepkowski**

Alexander Marquardt et al. proposed that the currently increased reality shows actually have an extremely restricted field of view contrasted with the human vision. To restrict out-of-view protests, scientists have overwhelmingly investigated visual direction ways to deal with imagine data in the restricted (in-view) screen space. In this paper, we look at a creative non-visual direction approach in light of sound material signals with the best-in-class visual direction procedure EyeSee360 for restricting carefully concealed objects in expanded reality shows with restricted field of view. In client study, they have assessed both direction techniques concerning search execution and circumstance mindfulness. Even all the more in this way, the sound material technique gives a huge improvement experiencing the same thing mindfulness contrasted with the visual methodology. By broadening the headband with more vibration engines, a higher goal and hence a rising exactness would be conceivable. Utilizing an alternate engine driver innovation like direct thunderous actuators or piezo electrics would likewise be sensible as far as usable data transfer capacity and speed increase qualities to further develop precision and execution.

2.9 Title: Joint cross-modal and unimodal Features for RGB-D Salient Object Detection,2021**Authors: Niachang Huang, Yi liu**

Niachang Huang et al. proposed the RGB-D remarkable item location that is one of the fundamental assignments in PC vision. Most existing models center around examining productive approaches to melding the correlative data from RGB what's more, profundity pictures for better saliency discovery. In any case, for some genuine cases, where one of the information pictures has poor visual quality, combining cross-modular elements does not assist with further developing the recognition exactness, when contrasted with utilizing unimodal elements as it were. Considering this, an original RGB-D striking article location model is proposed by all the while taking advantage of the cross-modal highlights from the RGB-D pictures and the unimodal elements from the information RGB and profundity pictures for saliency discovery. A Feature Selection Module is planned to adaptively choose those exceptionally discriminative highlights for the last saliency forecast from the melded cross-modular elements and the unimodal elements. A clever RGBD notable article discovery model has been proposed in this paper, where the cross-level and cross-modular highlights from the RGB-D picture matches, and the profundity pictures, are all the while caught and saved during the combination interaction by utilizing a proposed MFFM.

2.10 Title: Towards a New Learning Experience through a Mobile Application with Augmented Reality,2021**Authors: Santiago criollo, David Abad Vasquez**

Santiago et al. proposed a new learning experience with the ascent of data innovation and digitization, training has been confronted with the need to embrace new learning models utilizing innovation to make creative instructive approaches. Likewise, because of pandemic limitations and to assist with containing the spread of the infection (COVID19), all instructive organizations have been compelled to switch promptly to online training. The utilization of expanded reality (AR) in schooling gives significant advantages, expanded commitment and intuitiveness, and can assist with limiting the adverse consequences of the interruption of up close and personal schooling. Accordingly, this paper centers around depicting the impact of an expanded reality versatile application (NetAR) that was produced for designing understudies as a supplement to customary schooling. In this way, the convenience of the application was assessed with the IBM Computer System Usability Survey (CSUQ) device.

III. PROPOSED SYSTEM

The proposed systems consist of various stages, (i) Pre-processing, (ii) Edge Detection, (iii) Screen pointing, (iv) Homogeneous Device Interaction, (v) Performance measure. The architecture of the proposed systems is shown in Figure 1.

BLOCK DIAGRAM

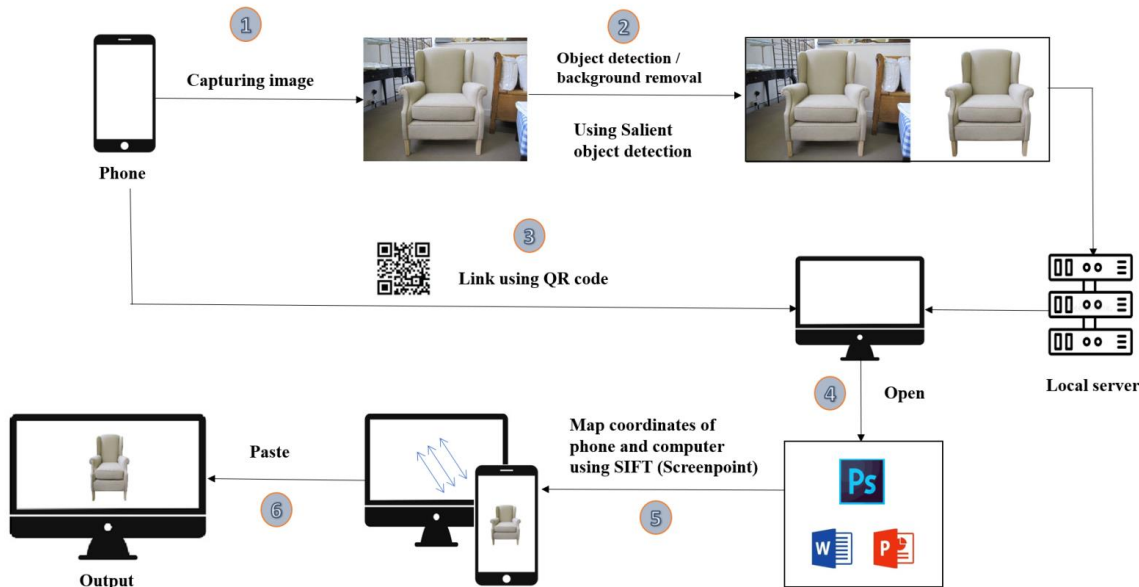


Figure 1 Proposed Object Detection System

3.1 PRE-PROCESSING

The point of pre-processing is to work on the nature of the image so the user can dissect it in a superior manner. By pre-processing one can smother undesired contortions and improve a few highlights which are important for the specific application the user is working for. Pre-processing is handling tasks with images at the least degree of abstraction both input and result are intensity images. These iconic pictures are of similar kind as the first information caught by the camera, with an intensity picture typically addressed by a matrix of image function values (brightness's). Albeit mathematical changes of pictures (for example turn, scaling, interpretation) are arranged among pre-handling strategies here since comparable methods are utilized. Here, the input image that has undergone pre-processing is further fed into the next module i.e., edge detection.

3.2 EDGE DETECTION

Edge recognition is an image handling method for tracking down the limits of objects inside pictures. It works by identifying discontinuities in lambency. Edge detection is utilized for image segmentation and information extraction in regions, for example, image handling, computer vision, and machine vision. In a picture, an edge is a curve that follows a way of quick change in picture intensity. Edges are frequently connected with the limits of items in a scene. Edge location is utilized to recognize the edges in a picture. To observe edges, the edge function can be utilized. This capacity searches for places in the picture where the power changes quickly, utilizing one of these two standards: Where the primary subordinate of the power is bigger in extent than some edge Where the second subsidiary of the power has a zero intersection Edge gives a few subsidiary estimators, every one of which carries out one of these definitions. For a portion of these estimators, it can be determined whether the activity ought to be touchy to even edges, vertical edges, or both. Edge returns a paired picture containing 1's in the same place as edges found and 0's somewhere else.

3.3 SCREEN POINTING

Screen point is characterized in pixels. The base left of the screen is (0,0); base right of the screen is (pixelWidth,0), left-top is (0, pixelHeight)

then, the right-top is (pixelWidth, pixelHeight). How about if this is comprehended with a picture; in this way, assuming the screen resolution is 720×480. The base left of the screen is (0,0); base right of the screen is (720,0), upper left is (0, 480) and the upper right is (720, 480).

The middle point will be (360, 240). This point can be utilized to observe which a big part of the screen is contacted. On the off chance that

- MousePosition < pixelWidth /2, the touch is on the left half of the screen
- MousePosition >= pixelWidth /2, the touch is on the right half of the screen

3.4 HOMOGENOUS DEVICE INTERACTION

In the proposed system an app is created for smart phone (mobile) which consist of a scanner. It scans the QR Code which will be generated when running the local server in the computer. QR Code is a machine-readable code consisting of a matrix barcode of black and white squares, typically used to store URLs or other information that can be read by a smart

phone's camera. Once the scanning is done the mobile and computer will be connected so that the pasting of object can be done.

3.5 U² NET ALGORITHM

LEFT PATH: Contraction path

| | |
|--------|---|
| Step 1 | : Input image is given (e.g., Image dimension-572 x 572 x 1) |
| Step 2 | : The input goes through convolution twice with 64 channel |
| Step 3 | : The number of channel will change from 1 → 64, as convolution process proceeds |
| Step 4 | : The last layer of each block has 2 x 2 max pooling layer which half down the size of image (568 x 568 to 284 x 284) |
| Step 5 | : This process is repeated 3 times |
| Step 6 | : The bottom layer with no maxpooling is reached |
| Step 7 | : The image at this moment has been resized to 28 x 28 x 1024 |

Figure 2 U² NET Algorithm – LEFT PATH

RIGHT PATH: Expansive path

| | |
|--------|---|
| Step 1 | : Transposed convolution (unsamplinnng technique) process is eventuated |
| Step 2 | : Then the image is upsized from 28 x 28 x 1024 -> 56 x 56 x 512 |
| Step 3 | : Then this image is concatenated with the corresponding image from the contracting path and together makes an image of size 56 x 56 x 1024 |
| Step 4 | : The processed image is then passed through multiple convolution each of 4 x 4 |
| Step 5 | : This whole process is repeated 3 times |
| Step 6 | : As the last step, the reshaping of image to satisfy our prediction requirements |
| Step 7 | : The processed image is the desired output |

Fig. 3 U² NET Algorithm RIGHT PATH

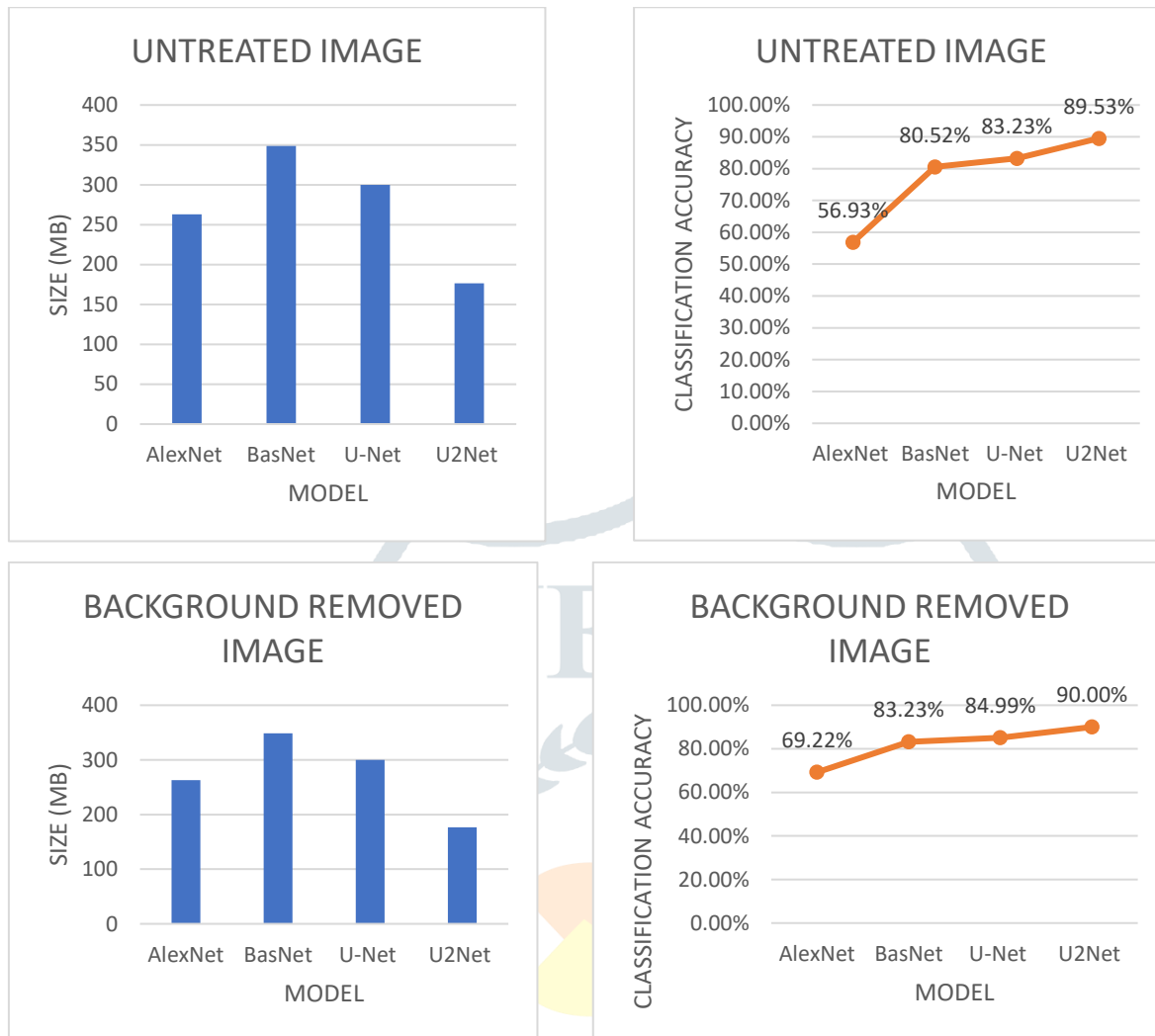
3.6 PERFORMANCE MEASURE

i) Classification Accuracy, $X = \frac{t}{n} * 100$

t – Number of correct classifications

n – Total number of samples

| | MODEL | SIZE(MB) | TOTAL NO. | CLASSIFICATION ACCURACY |
|--------------------------|-----------------------------------|-----------------|------------|-------------------------|
| Untreated image | AlexNet (Alex Krizhevsky et al) | 263.0MB | 10K | 56.93% |
| | BasNet (Xuebin Qin et all) | 348.5 MB | 10K | 80.52% |
| | U-Net (Olaf Ronnerberger et. al) | 300.0 MB | 10K | 83.23% |
| | U²Net | 176.5 MB | 10K | 89.53% |
| Background Removed Image | AlexNet (Alex Krizhevsky et al) | 263.0MB | 10K | 69.22% |
| | BasNet (Xuebin Qin et all) | 348.5 MB | 10K | 83.23% |
| | U-Net (Olaf Ronnerberger) | 300.0 MB | 10K | 84.99% |
| | U²Net | 176.5 MB | 10K | 90.00% |



$$\text{ii) Accuracy} = \frac{(tp+tn)}{(tp+tn+fp+fn)}$$

$$\text{Sensitivity} = \frac{tp}{tp+fn}$$

$$\text{Specificity} = \frac{tn}{tn+fp}$$

were, **tp** – true positive

fp – false positivity

tn – true negativity

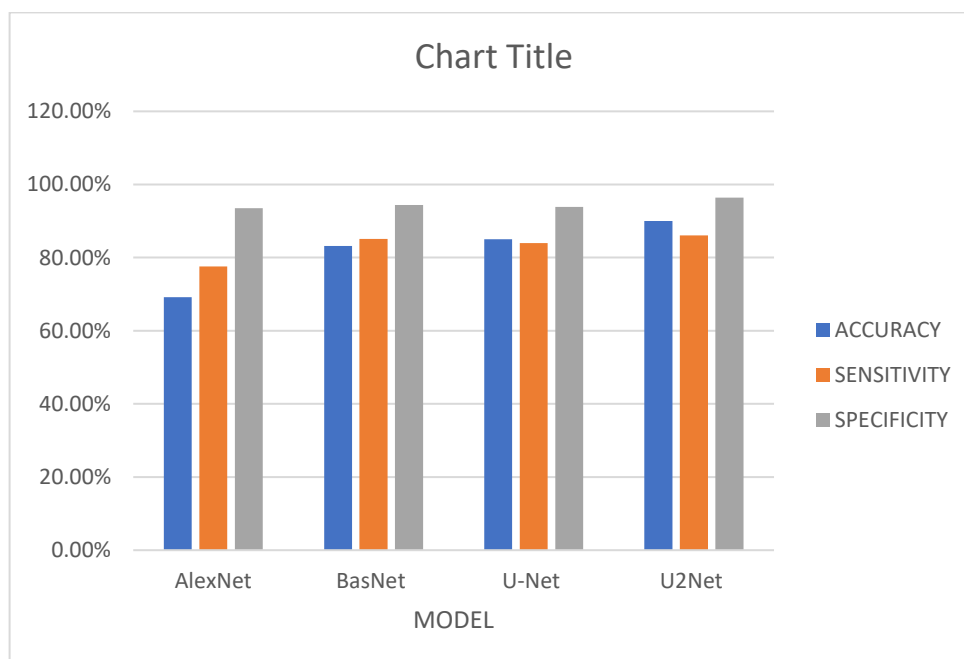
fn – false negativity

IV. RESULTS AND EXPERIMENTAL DATA

To evaluate how the model trained on U² Net generalizes, we use the comparison between AlexNet, BASNet and U-Net. We capture an object from an image using mobile phone camera within the field of view. The phone is the bought in front of the computer screen in order to map the coordinated of computer and phone via Screen point with comes under Scale Invariant Feature Transform. Once the coordinates are mapped, the transfer of the edge detected object from mobile phone to computer is done by connecting both the computer and mobile phone with same internet.

To evaluate with prior work, we provide the comparison of BASNet and U² Net on table 2. To minimize the time consumption and increase the accuracy, sensitivity, selectivity the new U² net algorithm is used. In addition, ClearGrasp algorithm is included to capture the transparent objects in an image. Our method significantly outperforms existing methods and demonstrations.

| MODEL | ACCURACY | SENSITIVITY | SPECIFICITY |
|----------------------------------|------------|---------------|---------------|
| AlexNet (Alex Krizhevsky et al) | 69.22% | 77.6% | 93.5% |
| BasNet (Xuebin Qin et al) | 83.23% | 85.1% | 94.4% |
| U-Net (Olaf Ronnerberger et. al) | 84.99% | 84.0% | 93.87% |
| U ² Net | 90% | 86.07% | 96.43% |



V. CONCLUSION

The crop and drop tool let you to digitalize the real world objects around us. It is an innovative idea that allows the user to take photo and detect the object in the real world and drop the image into the desktop computer. This tool will make the job of users much easier, as they will be able to point their phone cameras at an object and copy-paste it on their computers, instead of taking a photo, editing it and then inserting the cut-out into the document. The tool uses Augmented Reality (AR) and machine learning algorithm to detect the objects and isolate the image so that the background is automatically removed. For detection of image and for removal the background an open-source technology called U² Net is used and a computer vision algorithm called as Scale-invariant feature transform (SIFT) matches coordinates on the phone with the computer screen allowing you to place digital captures in specific locations on your computer screen. As a whole copy and paste your surroundings' using AR is the quickest way to capture, extract and transfer anything around you. This tool reverses the process and brings physical things into the digital world.

VI. FUTURE WORKS

This project will be very useful in terms of making presentations, editing, and documents that require images. Most of the time while making a presentation we go to Google for images and spend a lot of time finding the image, sometimes when we find it, there may be problems with pixels or quality so that we have to edit those images to get our desired result. This application will save a lot of time as you can directly pick the image from your surroundings and pass it

VII. REFERENCES

- [1] Chenggang yan, Xiaofei Zhou . “Edge Aware Multiscale Feature Integration Network for Salient Object Detection in Optical Remote Sensing Images”, IEEE Transaction On Geoscience and Remote Sensing, 2022, Vol 60, pp. 1-15.
- [2] Xuebin Qin, Zichen Zhang, Chenyang Huang, Chao Gao, Masood Dehghan and Martin Jagersand - “BASNet: Boundry-Aware Salient Object Detection”, IEEE Xplore, 2022, pp. 1-11.
- [3] Hao Wei, Rui Chen - “BASNet :A Boundry Aware siamese Network for accurate remote sensing change detection”, IEEE Geoscience and Remote Sensing Letters, 2022, Vol 19, pp. 1-5.
- [4] Xuebin Qin, Zichen Zhang, Chenyang Huang - “U²-Net going deeper with nested U-structured for Salient Object detection”, Elsevier, 2022, pp. 1-15.
- [5] Min Qiao, Gang Zhou, Qiu Ling Liu and Li Zhang - “Salient Object Detection: An Accurate and Efficient Method for Complex Shape Objects” , IEEE Access, 2021, Vol 9, pp.1-11.
- [6] Qiudan Zhang, Shiqi Wang - “A multi-task collaborative network for light field salient object detection”, IEEE Transactions on circuits and systems for video technology, 2021, vol 31, Issue no: 1849, pp.1-13.
- [7] Nianchang Huang, Yi Liu - “Joint Cross-modal and Unimodal Features for RGB-DSalient Object Detection”, IEEE Transactions on multimedia, 2021, Vol 23, Issue no:2428, pp. 1-14.
- [8] Santiago criollo, David Abad Vasquez - “Towards a New Learning Experience through a Mobile Application with Augmented Reality in Engineering Education”, Journal of MDPI, 2021, Issue no: 4921, pp.1-18.
- [9] Bing Zhou, Sinem - “Fine-Grained Visual Recognition in mobile Augmented Reality for Technical Support”, IEEE Transactions on visualization and computer graphics, 2020, Vol 26, Issue no:3514, pp. 1-10.
- [10] Alexander Marquardt, Christina Trepkowski – “Comparing Non visual and visual Guidance Methods for Narrow field of view Augmented reality Displays”, IEEE Transactions on visualization and computer graphics, 2020, Vol 26, Issue no:3389, pp. 1-13.