# A Machine Learning and Augmented Reality based Framework for Multilingual Product Identification in Retail using Mobilenets and Vuforia

**Geetanjali Bhola, Amogh Bansal, Divij Aggarwal, Gagan Kishor Upadhyay**

*Abstract*: *In the modern world, Machine Learning and Augmented Reality have taken the retail industry by storm. Machine Learning and Augmented Reality have provided a major boost to the industry of interactive retail by providing features such as real-time product detection and identification. The proposed research aims at overcoming several challenges in the present scenario which include the time consuming process of standing in long queues while purchasing the products at supermarkets, personalizing the shopping experience in order to maintain the privacy of the users, helping the customers to maintain their specified budgets, reducing the high labor costs and overcoming the language barriers while pertaining to selling products to groups with different linguistic backgrounds by combining the real-world interaction of Augmented Reality and Machine Learning-based product identification. This proposed research work aims at providing the customers with a futuristic shopping experience while maintaining their specified budgets. The Machine Learning-based object detection approach detected the products with around 96% accuracy and the Vuforia-based Augmented Reality approach detected objects with maximum accuracy.*

*Keywords*: *Augmented Reality, Vuforia, Machine Learning, MobileNets, Single Shot Multi-Box Detector.*

## I. INTRODUCTION

Augmented Reality and Machine Learning have created quite a buzz in the retail industry. Retail in today's economy can be combined with the functionalities of object detection and Augmented Reality to improve the customer experience.

Nowadays, customers tend to use high-end retail services but the retail industry presently faces a number of obstacles such as reducing the labor costs and eliminating the process of standing in long queues while purchasing various products. One of the major challenges faced by the retail industry is the language barrier across different regions while maintaining budget constraints which makes shopping across borders quite a tedious task. This research work aims at providing retail services using the real-world interaction of Marker-less Augmented Reality and Machine Learning-based object detection. This provides users with a personalized multi-lingual shopping functionality so as to maintain the privacy of the users [1] with features such as accurate real-time product detection, display of nutritional contents, product characteristics and video advertisements while maintaining the specified budget constraints, presentation of purchase history using login credentials and SQL databases, all available in multiple languages across different regions. The ability to allow the customers to view the cart and instructions on how to use the product in their local language improves the overall shopping experience of the customer [2]. Object detection helps to locate certain classes in the image. Augmented Reality allows coupling and superimposition of computer-generated imagery with real-world entities [3]. AR also helps to improve the user interface (UI) of the applications, improve the speed of local target recognition, use of technology that provides highly robust target-tracking of low-light and partly covered targets while pertaining to language barriers [4]

## II. CONCEPTUAL FRAMEWORK

### A. Object detection using TensorFlow

The first step involved use of TensorFlow for Object Detection of products in case of the Machine Learning based approach. The TensorFlow Object Detection API is an open source framework which simplifies the detection of objects and makes construction, training and deployment of object detection models easy as it uses pretrained object detection models. It also acts as a medium for detecting objects and algorithmic implementation. TensorFlow is given the name because it takes multidimensional arrays, also known as tensors as input to the model. It can be used with small changes on the several diverse frameworks.

**Ms. Geetanjali Bhola\*,** Department of Information Technology, Delhi Technological University, New Delhi, India. Email: geetanjali.getz@gmail.com
**Amogh Bansal,** Department of Information Technology, Delhi Technological University, New Delhi, India. Email: amoghdtu@gmail.com
**Divij Aggarwal,** Department of Information Technology, Delhi Technological University, New Delhi, India. Email: divij.aggarwal1410@gmail.com
**Gagan Kishor Upadhyay,** Department of Information Technology, Delhi Technological University, New Delhi, India. Email: gagankishor98@gmail.com

The framework is flexible in the sense that it is able to adjust to new conditions and hence, can be used for expressing several different algorithms.

### B. Object detection using MobileNets and SSD

MobileNet [5] is based on depth wise convolution in the case of which each filter contains a single input. A light MobileNet based model must include at least five layers with feature size 14×14×512. There can be 2 types of MobileNets design models, thinner or shallower in which the thinner model is more efficient than the shallow model.

Before the SSD, several different algorithms for object detection like R-CNN, Fast- RCNN etc. were obtained by us. R-CNN [6] works by suggesting different spots present in an image that act as input to an object classifier at the time of training. Fast- RCNN [7] and Faster-RCNN [8] are the improved versions of R-CNN, they use some techniques so as to reduce the manipulation of proposed regions. The problem of inference speed for real life objects was overcome by the YOLO [9] system. YOLO was not able to identify smaller objects in the frame, which was its main limitation. The problem was solved by SSD [10], the mAP value approached 74 percent and detection speed was similar to that of YOLO [11].

SSD serves the purpose of real time detection of various objects. Faster R-CNN makes use of a region proposal network which creates boundary boxes so that the objects belonging to several different classes can be detected whereas in case of SSD, region proposal network is no longer needed as it helps to fasten the process and hence it is eliminated. To regain the accuracy, SSD makes minor technical advancements including the default boxes and multi-scale features, using which, SSD is capable of meeting the accuracy achieved by Faster R-CNN with the use of low-resolution object images. The working of SSD was divided in two parts i.e. extracting the feature maps and application of convolution filters. SSD makes the use of small convolution filters to compute both the class and location scores. After the above steps, in order to make predictions about the classes of the detected objects, feature maps are extracted and 3×3 convolution filters are applied on all cells individually. MobileNets and SSD provide the benefit of achieving lag-free accurate object detection on a system with nominal specifications.

### C. Creation of SQLite Database in Machine Learning

The next step involved the creation of SQLite Database which serves as an open-source SQL database engine that doesn't require a separate server to operate. It is an in-process library that implements a self-contained, zero-configuration, serverless, transactional SQL database Engine that has bindings to several programming languages such as Java, Python, C++, Basic [12]. In our case we have made use of Python so as to create a database in case of the Machine Learning based approach which serves the purpose of storing the login credentials as well as the purchase history pertaining to the previous visits of the users.

### D. Vuforia based Image Recognition and superimposition of virtual information

In order to detect products in Augmented Reality, we make use of Vuforia which is an AR software development kit which was developed by Qualcomm and allows the creation of various Augmented Reality applications. It allows the recognition and tracking of planar images and 3-D objects in real-time. Virtual objects are placed through a camera viewfinder and then the position of objects is adjusted in the background of the objects and in front of the lens.

In the development process, Vuforia uses target images but any kind of image does not qualify as a target image. All the uploaded images, at first, are evaluated by the Vuforia web developer. Scoring of value target defines how well an image can be detected and tracked using Vuforia SDK. The score of an image acts as a target range of 0 to 5 which helps in assessment [13]. Greater number of stars allows to obtain stronger detection and tracking capabilities of an image target.

In order to determine the number of stars, various features serve as the attributes of the assessment criteria. Repetitive patterns should not be there in the image target and the image target should contain high local contrast, uniform distribution of features and high feature density, so as to obtain better quality of recognition of the image.
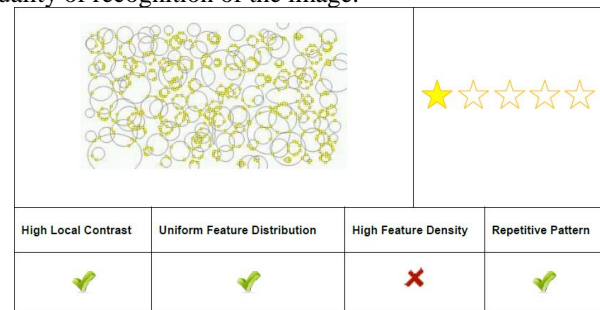


**Fig. 1. IT with 1-star ball shape**

For example, Fig. 1 which has been shown above, has a quality level in one star because the density of the features is poor. Hence, alteration of its line shape or texture is required, so as to increase (multiply) the number of features of this figure [14].



**Fig. 2. IT with 4 stars square shape**

Now, Fig. 2, which has been shown above, depicts the rise of quality in star level of the image target by increasing the number of features to improve feature density and similarly for different image targets the quality in star level can be improved by also increasing the local contrast, distributing the features uniformly and through avoidance of repetitive patterns.

This information was integrated in the research work by implementing Marker-Less Augmented Reality which helped in obtaining stronger detection and tracking capabilities of the image targets.

The ability of augmented reality technology to blend real scenes and virtual information was used for the superimposition of nutrient distribution and advertisement over a particular product along with information about the remaining budget and the shopping cart.

### E. Database Creation using MAMP in AR

Creation of a database in the Augmented Reality approach was carried out using MAMP (Macintosh, Apache, MySQL, and PHP) which allows you to have access to a local PHP server and MySQL server and provides you with all the tools necessary to run WordPress on your machine, for testing and development purposes. The created database allowed storage of the login credentials of the different users in the Augmented Reality-based approach.

### III. METHODOLOGY

The main aim of the research was to provide consumers with an interactive experience through real-time product identification in the retail sector.

The flow of the implementation of this idea has been shown in Fig. 3. At first, object detection was carried out with the help of TensorFlow, which because of its high GPU intensiveness and low efficiency resulted in low accuracy and fps, along with lag in product detection. So as to overcome these problems, product detection was therefore implemented using a joint version of Single Shot MultiBox Detector (SSD) structure and Mobile Net classifier [15].

On observation of improvement in detection accuracy and decrease in lag in the latter method, new features were added to enhance the overall user experience. Features such as advertisement, description of the products and interactive budget management for customers were added. Then, a database was created using Python so as to store the login credentials of various users and also provisions for different languages were provided in case of the Machine Learning-based approach. Multiple windows were being required for different features which posed to be a limitation of the Machine Learning-based approach, hence the focus was shifted towards incorporation of Augmented Reality in the research.

At first, the implementation of marker-based AR was done, which performed real-time superimposition of nutritional contents on their respective markers for various products. The marker-less technique was executed due to the additional cost of production of markers, and using marker-less technique, real-time superimposition was done without the use of predefined markers for various products.

Then in order to provide a better customer experience, new features such as display of advertisements, nutrition pie charts of different products as well as interactive budget management with shopping-mart guide for users. Finally, a database was created using MAMP so as to store the login credentials of the users and the option of choosing between various languages was provided.
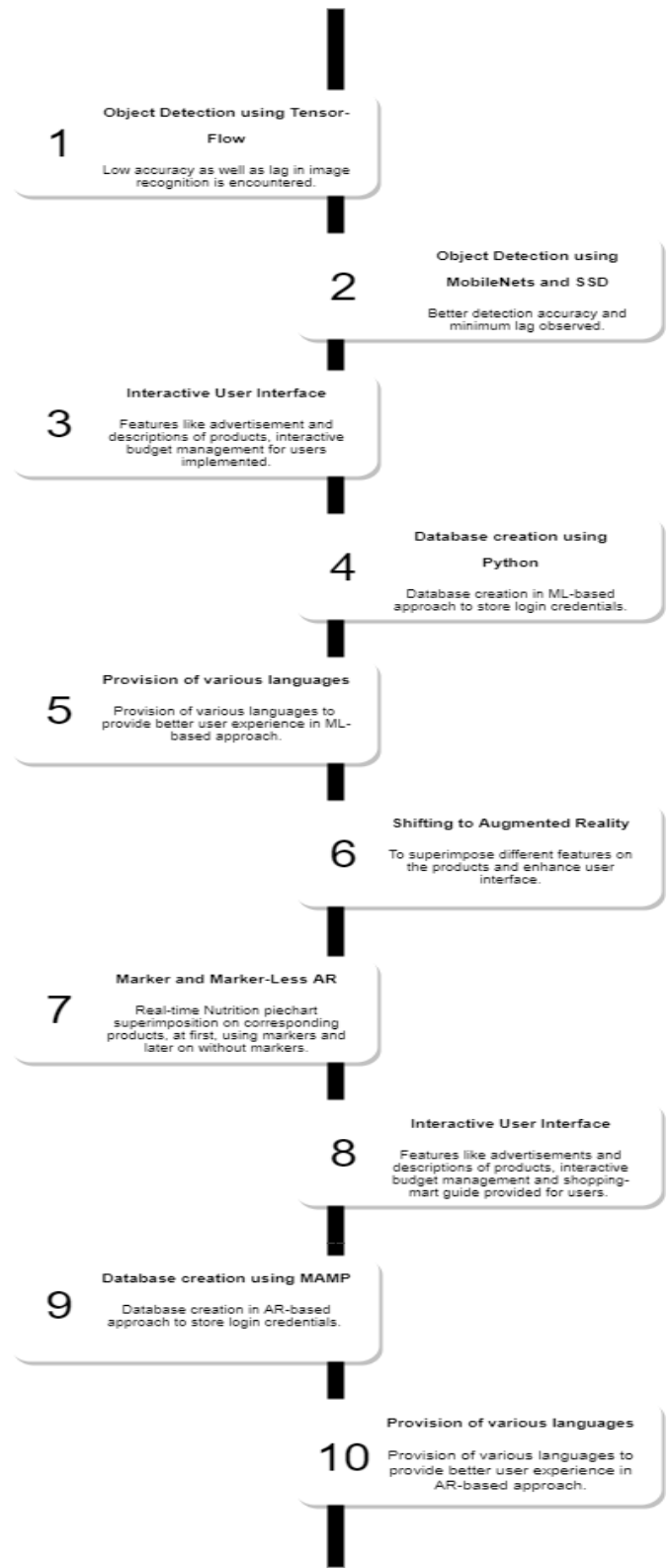


**Fig. 3. Project Flow**

### A. Technology Stack

#### 1. OpenCV

OpenCV is a large open-source library for Machine Learning, Computer Vision as well as Image Processing. It plays a great role in real-time operations. It performed the task of processing images and videos to detect objects and products.

1924

OpenCV was the first piece of technology used to detect objects using Python in case of the Machine-Learning based approach.

### 2. AR Camera

In the case of Augmented Reality, AR Camera was needed to be imported in unity which is used to provide AR scenes. It takes the help of WebCamBehaviour class, VuforiaBehaviour class, DatabaseLoadBehaviour class to accomplish its functions. The WebCamBehaviour class, which also serves as the prime camera control class, informs about whether to switch on the main camera, which camera to use etc. The control entry of Vuforia is the VuforiaBehaviour class, which registers all the open interfaces of Vuforia. The DatabaseLoadBehaviour class is used to trigger identification library, which can perform the function of loading multiple databases that have been uploaded to Vuforia for processing and specifying which library to trigger [13].

### B. Implementation

### 1. Coco Dataset for mobilenets and SSD based object detection

Common Objects in Context is a large-scale object detection, segmentation, and captioning dataset. COCO includes several features such as object segmentation, recognition in context, 330K images and 1.5 million object instances.

We have created a dataset similar to COCO dataset [16], for performing the Machine Learning based Product Identification, that has images of shopping objects like bottles and chairs. Nearly 300 images were taken of both the objects, amongst those images, 60% were utilized in the training phase and 40% in the testing phase. On the detection of multiple objects by the web camera, then a bounding box is drawn over every single detected object with its class as the heading of the bounding box.

### 2. Database Creation for image targets (Quality of Image targets)

The next step was the construction of the database of image targets for Augmented Reality based Object Detection using Vuforia SDK following which the license key of the created database was added to the research work. At first, the database consisted of the stickers which were used as image targets for products in marker-based object detection in AR. Since the initial rating of the stickers was low (as mentioned in Table 1), hence weaker detection and tracking ability of the image were experienced. Then so as to eliminate the concept of stickers and enhance the trackability of the image, Marker-Less product identification was implemented in AR. The database that was prepared for the marker-less execution consisted of images of actual products which were used as image targets for superimposition. The rating of both the image targets improved significantly (Refer Table 1) and reached 5 stars that ensured better trackability and detection of the images and hence enhanced the functioning of the model.

**Table I- Target Score**

| No. | Target Image (Marker-Based) | Star | No. | Target Image (Marker-Less) | Star |
|-----|------------------------------|------|-----|-----------------------------|------|
| 1 | Apple Sticker | 2 | 1 | Knorr Soup Image | 5 |
| 2 | Maggie Sticker | 2 | 2 | ACT II Popcorn Image | 5 |
| | *Average (Marker-Based)* | *2* | | *Average (Marker-Less)* | *5* |

### 3. Provision of Optional Languages

The option of choosing from various languages was then provided in both the Machine Learning and Augmented Reality based approaches in order to overcome the challenges of language barriers across different regions as per the requirement of the users.

### 4. Database Creation for purchase history and Login Credentials

Finally, databases were created for storing the purchase history as well as the login credentials of the users. In case of the Machine Learning-based approach the SQLite database [12] was created using Python and in case of the Augmented Reality-based approach a database was created using MAMP so as to store the login credentials and the purchase history of the users.

## IV. RESULTS

We have tested the machine learning part of our research on an i5 processor laptop. In the machine learning based implementation of the research work, at first, the object detection is carried out using TensorFlow (Fig. 4), but due low accuracy as well as lag in image recognition, the focus was shifted towards object detection using MobileNets and SSD (Fig. 5). After that, the user is made to login (Fig. 6) which provides the user an option to view the user's purchase history (Fig. 7) and then the user can choose a language of his/her choice (Fig. 8). Then, the user enters the budget for shopping (Fig. 9) and after the budget is provided, the process of object detection starts. Once the object is detected successfully (Fig. 10), then the user is asked whether he/she is interested to view the description of the product (Fig. 11) or viewing the advertisement associated with the product. After the user makes the choice, the user has the option to either buy it or ignore it. If the user decides to buy the product, then the item is added to the cart. The user can access the cart anytime to know about the existing products in the cart and the total amount. However, when a user exceeds the budget entered by the user initially, a pop-up showing OVER-BUDGET notification appears and further purchasing of items can't be performed (Fig. 12). Hence, this whole system allows the user to automatically detect any product of interest and then take the decision of buying it or not after seeing the advertisement and nutritional contents of the product.
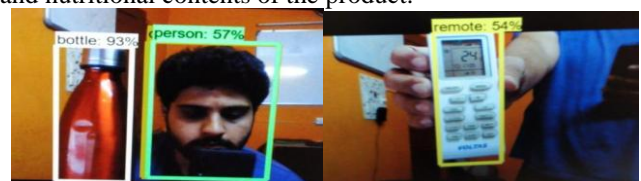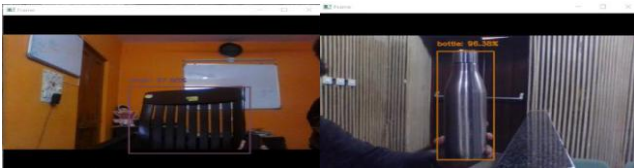
**Fig. 4. Object detection using TensorFlow**



**Fig. 5. Object detection using MobileNets and SSD**



**Fig. 6. User is logged in**



**Fig. 7. User can view Purchase History**
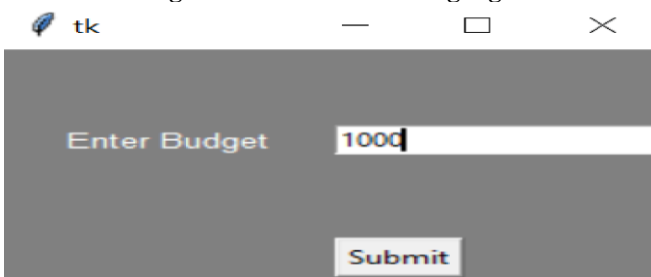


**Fig. 8. User can choose language**



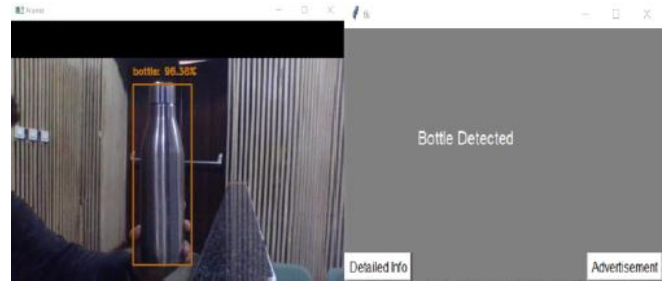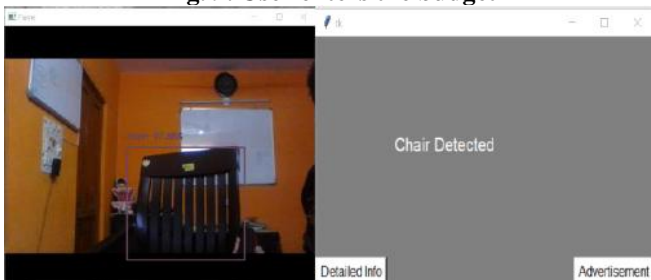**Fig. 9. User enters the budget**

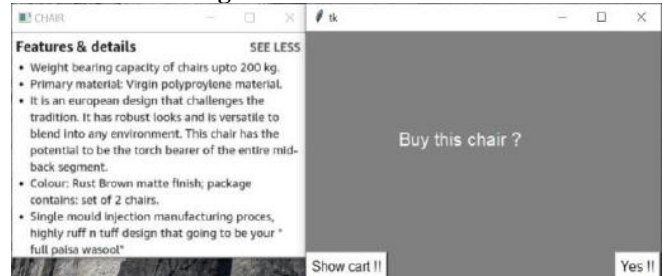



**Fig. 10. Product detected**



**Fig. 11. User can view description**



**Fig. 12. Over Budget**

The main limitation of this approach was that there were several pop-up windows through which the user had to go through to reach the final step, and hence, we then decided to perform the same implementation in Augmented Reality, where the user could interact with the features easily without switching between several windows.

In the marker-less AR implementation of the research work, initially the user is not registered (Fig. 13). Once the user completes the registration process (Fig. 14), his/her login credentials are visible in the database (Fig. 15). After that the user can sign into the login portal using his/her login credentials (Fig. 16). Once the user was logged in (Fig. 17), the product was scanned by the Vuforia SDK and the UI consisting of panels and shopping-mart guide was superimposed on the product (Fig. 18). After that, the user is made to enter the budget for shopping which is then displayed on the panel superimposed at the top of the product. The panel which is superimposed on the left side of the product shows the various choices for the users which include choosing the language preferred by the user, display of advertisement and nutrition value of the product along with add to cart and show cart options (Fig. 19). After viewing the advertisement and the nutrient value of the product, the user has the option to buy the product or not. If the user chooses to buy it, then the product is added to the cart and the budget is updated. The user has also been provided with an option to view the cart and the remaining budget (Fig. 20). Finally, the user can also view his/her purchase history till date (Fig. 21). Hence, this whole model provides an opportunity for users to select the products which are in their budget as well as are perfectly suited for their requirements.

**Fig. 13. User is not registered**


**Fig. 14. Registration process**


**Fig. 15. Login credentials of the user are added to the database**


**Fig. 16. User logs in using his/her login credentials**


**Fig. 17. User has logged in**


**Fig. 18. Product has been scanned**


**Fig. 19. Product scan with language change, budget entry, language entry, video and nutrition pie-chart display**


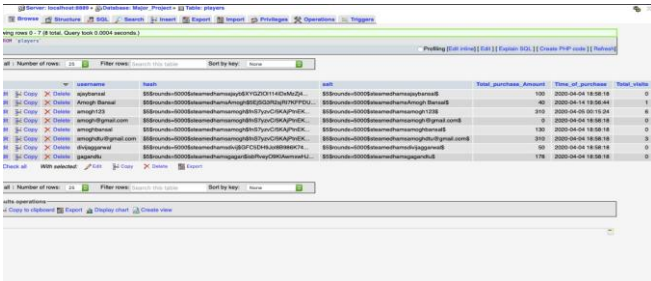**Fig. 20. Product scan with show cart and remaining budget**


**Fig. 21. Purchase history**

The summarized test cases of the important functions of the AR (Marker-Less) implementation of the research work is shown in Table 2.

**Table II- Test cases for important functions**

| No. | Precondition | Operation | Expected Results | Tested Results |
|-----|--------------|-----------|------------------|----------------|
| 1. | User not logged in | "Register button is clicked" | User is redirected to the registration portal | Pass |

| | | | | |
|---|---|---|---|---|
| 2. | Login Credentials of new users are added | "Register" button is clicked | Registration of new user complete, login credentials are added to database and the user is redirected to the login portal | Pass |
| 3. | User Registered | Login Credentials of existing users are added | Login Credentials are displayed in text format on the panel | Pass |
| 4. | Login Credentials of existing users are added | "Login" button is clicked | User is logged in | Pass |
| 5. | User is logged in | "Scan Product" button is clicked | Product is scanned | Pass |
| 6. | Product is scanned | Language of user's choice is entered | Changes displayed according to the language entered | Pass |
| 7. | Product is scanned | Click on the shopping mart guide | "Hi! Welcome to mart" text is displayed and waving animation is performed by the guide | Pass |
| 8. | Product is scanned | "Advertisement" button is clicked | Product advertisement with audio is shown in 3-D format | Pass |
| 9. | "Advertisement" button is clicked | "Advertisement" button is clicked again | 3-D advertisement is removed | Pass |
| 10. | Product is scanned | "Nutrition" button is clicked | Product's nutritional pie chart with audio is shown in 3-D format | Pass |
| 11. | "Nutrition" button is clicked | "Nutrition" button is clicked again | 3-D display of product's nutritional pie chart is removed | Pass |
| 12. | Budget is entered | Enter Budget for shopping | Budget is displayed in text format on the panel | Pass |
| 13. | Budget is entered | "Add to Cart" button is clicked | Product is added to the cart and the updated budget (Current budget – price of the product scanned) is displayed | Pass |
| 14. | Product is added to the cart | Wait for a few seconds | Shopping mart guide displays "Thank you for choosing (product name)" in text format on the panel | Pass |
| 15. | Products are added to the cart | "Show Cart" button is clicked | Products in the cart are displayed in text format on the panel | Pass |
| 16. | Budget becomes zero | "Add to Cart" button is clicked | "Over Budget" is displayed in text format on the panel | Pass |
| 17. | Budget becomes zero | "Add to Cart" button is clicked | Shopping mart guide displays "Please check your budget" in text format on the panel | Pass |
| 18. | Shopping mart guide is clicked | Shopping mart guide is clicked again | Text displayed by shopping mart guide is removed | Pass |
| 19. | Product is scanned | "Purchase History" button is clicked | The entire purchase history of the customer including the customer's name, total amount spent till date, date and time of last purchase and total number of visits till date are displayed in text format on the panel | Pass |
| 20. | At least 1 product is scanned | Scan another product | Features like budget and cart content are transferred to the newly scanned product's user interface (UI) | Pass |

## V. CONCLUSION AND FUTURE WORK

The research provides the retail functionality combined with the real-world interaction of Augmented Reality and Machine Learning-based product identification, thus providing a fresh experience to the customers. Supermarket shopping is provided in a stimulating and interesting manner to the customers.

Hence, the research work's design goals were achieved and therefore, research work's future vision includes increasing the scalability of the research work to include a larger database of users, provisions for more languages, a wide range of products in various industries such as sports, footwear, groceries and clothing Also, future work includes implementation of a recommender system for providing product suggestions to customers.

## REFERENCES

1. A. Poushneh, "Augmented reality in retail: A trade-off between user's control of access to personal information and augmentation quality," Journal of Retailing and Consumer Services, vol. 41, pp. 169-176, 2018.
2. A. Poushneh and A. Z. Vasquez-Parraga, "Discernible impact of augmented reality on retail customer's experience, satisfaction and willingness to buy," Journal of Retailing and Consumer Services, vol. 34, pp. 229-234, 2017.
3. J. Carmigniani and B. Furht, "Augmented reality: an overview," in Handbook of augmented reality, Springer, 2011, pp. 3-46.
4. F. Peng and J. Zhai, "A mobile augmented reality system for exhibition hall based on Vuforia," 2017.
5. A. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto and H. Adam, "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications," 2017.
6. R. Girshick, J. Donahue, T. Darrell and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," 2014.
7. R. Girshick, "Fast r-cnn," in Proceedings of the IEEE international conference on computer vision, 2015.
8. S. Ren, K. He, R. Girshick and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in Advances in neural information processing systems, 2015.
9. J. Redmon, S. Divvala, R. Girshick and A. Farhadi, "You Only Look Once: Unified, Real-Time Object Detection," 2016.
10. W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu and A. Berg, "SSD: Single Shot MultiBox Detector," 2016.
11. S. Kanimozhi, G. Gayathri and T. Mala, "Multiple Real-time object identification using Single shot Multi-Box detection," 2019.
12. S. T. Bhosale, T. Patil and P. Patil, "Sqlite: Light database system," International Journal of Computer Science and Mobile Computing, vol. 4, no. 4, p. 882, 2015.
13. D. Adrianto, M. Hidajat and V. Yesmaya, "Augmented reality using Vuforia for marketing residence," 2016.

14. A. Simonetti Ibañez and J. Paredes Figueras, "Vuforia v1. 5 SDK: Analysis and evaluation of capabilities," 2013.
15. Y. Wang, C. Wang and H. Zhang, "Combining a single shot multibox detector with transfer learning for ship detection using sentinel-1 SAR images," Remote sensing letters, vol. 9, no. 8, pp. 780-788, 2018.
16. T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár and C. Zitnick, "Microsoft COCO: Common Objects in Context," 2014.

## AUTHORS PROFILE

**Ms. Geetanjali Bhola** is an Assistant Professor in Information Technology Department of DTU. She has been taking subjects like database, computer graphics, programming etc. She has co-authored many research papers in the field of image processing. With a deterministic attitude she has helped and encouraged her students to take up challenging research problems which have later won awards at different platforms. Her current research includes augmented reality applications, image processing to name a few.

**Amogh Bansal** is pursuing B. Tech. in Information Technology from Delhi Technological University, Delhi, India. He has published a research paper titled 'Augmented Reality and Machine Learning based Product Identification in Retail using Vuforia and MobileNets' and presented at the International Conference on Inventive Computation Technologies (ICICT-2020) held at Coimbatore, India from 23-26th February, 2020. He is also the Corporate Head of the Society of Robotics of Delhi Technological University (SRDTU).

**Divij Aggarwal** is pursuing B. Tech. in Information Technology from Delhi Technological University, Delhi, India. He has published a research paper titled 'Augmented Reality and Machine Learning based Product Identification in Retail using Vuforia and MobileNets' and presented at the International Conference on Inventive Computation Technologies (ICICT-2020) held at Coimbatore, India from 23-26th February, 2020. He has successfully mentored a batch of 50+ students in Data Structures and Algorithms in Java.

**Gagan Kishor Upadhyay** is pursuing B. Tech. in Information Technology from Delhi Technological University, Delhi, India. He has published a research paper titled 'Augmented Reality and Machine Learning based Product Identification in Retail using Vuforia and MobileNets' and presented at the International Conference on Inventive Computation Technologies (ICICT-2020) held at Coimbatore, India from 23-26th February, 2020. He also has participated and won the I.S.G.W Appathon held in New Delhi, 2017.