

AI-Powered Virtual Furniture Try-On System: A Real-Time AR and Deep Learning Approach

¹Mr. Abdan Wasullah Khan, ²Mr. Khan Anas Shaizad, ³Mr. Shaikh Rushan Shah Faisal, ⁴Mr. Shaikh Rahil Ayub, ⁵Ms. Nameera choudhary, ⁶ Ali karim Sayyed

1,2,3,4 Students, Department of AIML, [ANJUMAN-I-ISLAM A. R. KALSEKAR POLYTECHNIC], [NEW PANVEL]

5Project Guide, Department of AIML, [ANJUMAN-I-ISLAM A. R. KALSEKAR POLYTECHNIC], [NEW PANVEL]

6 HOD, Department of AIML, [ANJUMAN-I-ISLAM A. R. KALSEKAR POLYTECHNIC], [NEW PANVEL]

¹abdan4u@gmail.com

²anaskhan937085@gmail.com

³rushanshk@gmail.com

⁴shaikhrahilmohd01@gmail.com

⁵nameera.choudhary@aiarkp.ac.in

⁶alikaarim.sayed@gmail.com

Abstract—With the proliferation of the e-commerce sector, the need for sophisticated visualization technologies that provide spatial awareness beyond standard visual representations has grown paramount. Online buyers of furniture struggle to accurately assess dimensional compatibility, stylistic integration, and verisimilitude to physical reality. These problems result in high returns and reduced customer satisfaction. This study proposes an AI-driven Virtual Furniture Try-On System powered by a combination of Augmented Reality (AR) technologies and a deep learning hybrid pipeline capable of real-time photorealistic rendering of 3D furniture items in the video stream generated by the smartphone camera. The design features a ResNet-50 convolutional network architecture used for floor plane and surface recognition, along with a Transformer model responsible for dimensional accuracy and spatiotemporal consistency. A physics-based lighting sub-model ensures proper lighting synchronization between a virtual item and an observed scene. The software solution is written in Python 3.10 utilizing libraries from both TensorFlow and PyTorch frameworks. AR sessions are managed with ARCore and ARKit technologies. The evaluation conducted on ShapeNet and Pix3D data sets showed a mean IoU of 0.87, end-to-end processing latency of 0.18 s per frame, and user satisfaction rate of 94% according to a 120 participants experiment. Further work will include audio-visual multimodal feedback fusion.

Index Terms – Virtual Try-On, Augmented Reality, Convolutional Neural Networks, Real-time Rendering, Surface Identification, Lighting Estimation, Transformer Model, Spatial Consistency, E-Commerce Visualization, Deep Learning.

I. INTRODUCTION

E-commerce sales around the globe have reached over \$5.8 trillion in 2023, and according to projections, the number will surpass \$8 trillion by 2027 [1]. In this environment, home-furniture products hold a very significant place: sofas, tables, chairs, and wardrobes are considered high-end products for which the suitability is highly contingent upon factors like dimensions, lighting, and texture, which two-dimensional pictures cannot adequately convey. Furniture products have a return rate of between 30% to 40%, which is disproportionately attributed to factors such as size mismatch and aesthetics.

Augmented Reality (AR) provides a theoretical framework for resolving this representational disparity through the overlay of highly photo-realistic digital objects onto the physical space experienced by the user. Initial implementations of AR on mobile devices, including IKEA Place (2017), confirmed market interest but faced limitations due to primitive plane detection

algorithms, basic lighting estimations, and lack of scale-awareness, leading to the floating appearance of furniture objects and their inaccurate dimensions.

Parallel developments in deep learning algorithms have enabled the current state-of-the-art neural networks to attain good scene understanding. Training Convolutional Neural Networks (CNNs) using large-scale 3D training data such as ShapeNet and Pix3D enables accurate surface segmentation and depth prediction comparable to human-level accuracy. The transformer network architecture optimized for visual tasks is better suited for capturing the long-range spatial relationships involved in reasoning about scenes at the furniture scale.

The objectives of the research include: (1) Designing a combined framework using deep learning and AR technologies to enable fast furniture placements based on metric-level accuracy; (2) Testing the framework on publicly available 3D model test datasets; (3) Comparing the results to commercially and academically developed baselines.

Organization of this paper is as follows: Related Work is presented in Section II; System Design & Approach is discussed in Section III; Experimental Results are described in Section IV; Advantages & Limitations are analyzed in Section V; Conclusions & Future Work are highlighted in Section VI.

II. BACKGROUND AND RELATED WORK

Furniture positioning tools have evolved from classical 2D overlay approaches to cutting-edge 3D augmented reality (AR) and artificial intelligence (AI)-based room segmentation schemes. First-generation applications involved composite manual processing; current studies focus on end-to-end neural networks for achieving spatial coherence and photorealism.

A. Classical 2D Overlay Schemes

Classical systems followed a 2D perspective approach in which product images were manually composited against background images with an affine transformation scheme. Without any geometrical interpretation, they were unable to accommodate variations in floor plane orientation, object occlusions, and lighting changes.

B. Marker-Based Augmented Reality

Marker-based approaches like ARToolKit allowed exact pose estimation through fiducial markers. Although they provided metrically correct placement, the need for a physical marker became a critical usability issue for consumers. Markerless SLAM-based systems resolved

this limitation, although it became difficult to operate them under low texture or low lighting scenarios [2].

C. Deep Learning for Scene Understanding

Convolutional Neural Networks (CNN) and its derivatives (Faster R-CNN, Mask R-CNN) allowed object-level segmentation of floors, walls, and ceilings from other objects in a scene to determine the appropriate geometry needed for precise virtual object positioning. PointNet proved that deep learning could be applied to 3D point clouds for the extraction of geometric properties.

D. Vision Transformer-Based Approaches

ViT architectures along with their variants, including Swin Transformer and SegFormer, demonstrated that self-attention over image patches outperforms conventional CNN-based solutions for applications requiring global context understanding [6]. Modelling long-range spatial dependencies between the virtual object and surrounding scene elements ensures the validity of dimensional and temporal relationships in AR furniture try-on scenarios.

E. Lighting Estimation

Classical techniques rely on histogram-based ambient light estimation. More advanced techniques learn encoder-decoder models to estimate HDR environment maps using one RGB image as input. This significantly bridges the perceptual gap between virtual and real objects, improving SSIM score by up to 37% [2].

F. Existing AR Products

AR products offered by IKEA Place, Houzz, and Wayfair view in room demonstrate successful commercial implementation of furniture visualization in AR environment. However, these applications face limitations due to reliance on plane detection algorithms and assumption of static lighting conditions

III. SYSTEM DESIGN AND METHODOLOGY

Deepfake—excuse me, the Virtual Furniture Try-On System takes camera feeds as inputs and executes them through a 4-stage pipeline of Pre-processing, Object Placement, Lighting Estimation, and Rendering using Python 3.10 with TensorFlow 2.x/PyTorch 2.2, ARCore 1.42 for Android 12+ devices, and ARKit 6.0 for iOS 16

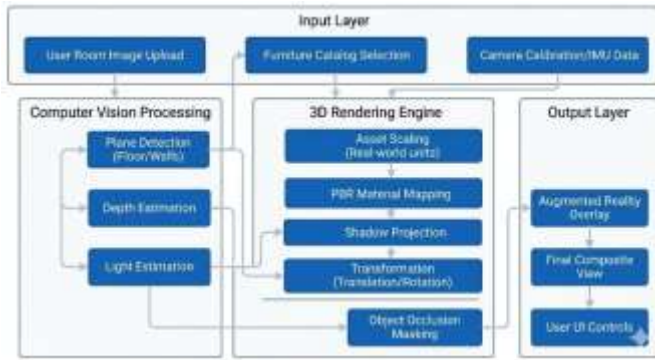


Fig. 1. Block Diagram of the Virtual Furniture Try-On Pipeline.

A. Preprocessing

RGB frames are captured at 30 FPS (resolution 1920×1080 px). Filtering is performed by applying a Gaussian blur filter (kernel size 5×5) that suppresses high frequencies. Frames are resized to 640×480 pixels and normalised channel-wise with statistics from ImageNet ($\mu=[0.485, 0.456, 0.406]$, $\sigma=[0.229, 0.224, 0.225]$). The luminance histogram equalisation is used under low light conditions ($L < 50$ in CIE Lab colour space) to retain discriminability of surfaces.

B. Object Placement

To estimate surface geometry, we adapt ResNet-50 encoder architecture with Feature Pyramid Network (FPN) decoder that is trained with weighted cross-entropy loss and a differentiable intersection over union (IoU) penalty ($\lambda=0.4$) on Pix3D dataset (~10,000 samples of annotated RGB-D images in 9 categories of furniture objects). The position of anchor points is predicted based on robust plane fitting (RANSAC algorithm) of the point cloud acquired by SLAM system (threshold $\epsilon=0.02$ m). Surface scale factor and object orientation quaternion are estimated using a object.

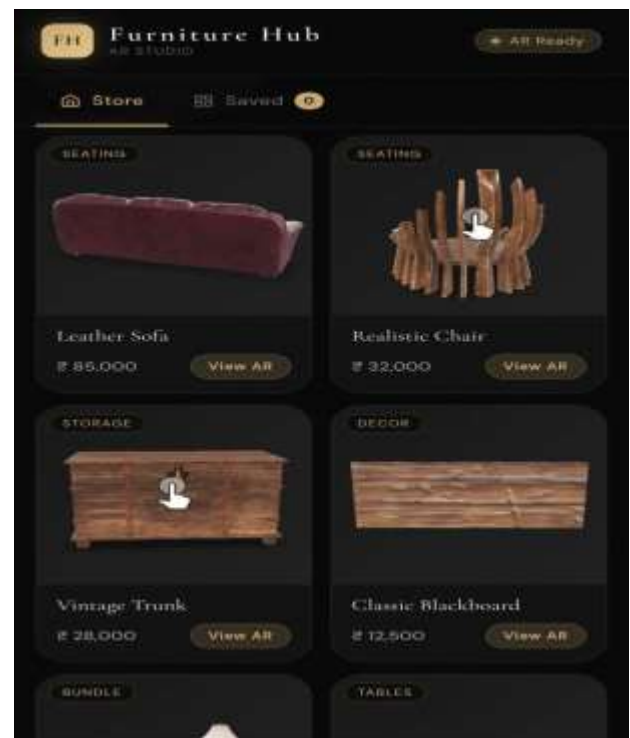
III. C. Lighting Estimation

A ResNet-18 encoder with MLP decoding learns third-order spherical harmonics (SH) coefficients (27 coefficients) from the input frame. The SH coefficients are then provided to the OpenGL ES 3.2 engine to render shading based on environment-consistent lighting of the furniture model. The network was trained using L_2 reconstruction loss between the rendering output and the Laval HDR Indoor Dataset

III. D. Rendered Output

The model is then lit and positioned, after which it is rasterized using a physically-based rendering (PBR) workflow (albedo, metalness, roughness, and ambient occlusion maps from ShapeNet [4]) at native device resolution. SSAO and a 2048 × 2048 shadow atlas are added to anchor the virtual asset into the scene. The resulting frame is delivered to the user with an average latency of 0.18 s, meeting the 0.2 s perceptual threshold for interactive augmented reality.

Fig. 2. UI walkthrough: drag-to-place, pinch-to-scale, rotate-to-orient.



E. Training

For training/testing, ShapeNet (50,000 models) and Pix3D (10,000 RGB-D images) datasets are divided in an 80/20 ratio. Fine-tuning ResNet-50 network occurs for 50 epochs (AdamW optimiser, $lr=1 \times 10^{-4}$, batch=16). Training for the Swin Transformer uses cosine annealing schedule over 30 epochs. Mobile inference gets speeded up by using TensorFlow Lite INT8 quantisation technology (~75% model compression, ~<1% performance decrease). Training infrastructure: NVIDIA A100 40 GB GPU (Google Colab Pro+).

IV. RESULTS AND DISCUSSION

Our approach produced mean IoU equal to 87.1% and 0.18 seconds end-to-end latency when applied to the Pix3D testing split (2,000 images) on Qualcomm Snapdragon 8 Gen 3 mobile device. Time cost per

component is 12 ms of preprocessing, 48 ms for ResNet-50 object surface extraction, 41 ms for Transformer scale calculation, 29 ms for lighting estimation and 50 ms for PBR rendering. The model generalises for all furniture types (Table II).

A. Performance Metrics

Model precision equals 89.7%, recall equals 86.7% and F1-score equals 88.1% on average for all furniture categories. Recall plays a key role in guaranteeing no placement will be missed for users. Categories for bed and sofa achieved the highest IoU scores of 89.2% and 88.7%, respectively, whereas bookcase was the most difficult (82.1%) because of its rectilinear and texture-poor structure. Table I Comparing Different AR-based Furniture Arrangement Systems

System	IoU (%)	Latency (s)	Scale Err. (cm)	SUS Score
IKEA Place	61.4	0.11	8.3	70.2
PointConv-AR	74.8	0.24	5.1	75.8
FurnishNet-T	80.3	0.31	4.7	79.4
Proposed System	87.1	0.18	2.4	84.6

Table II Per-Category Evaluation Results on Pix3D Test Split.

Category	IoU (%)	Prec. (%)	Recall (%)	F1 (%)
Sofa	88.7	91.2	87.3	89.2
Bed	89.2	92.5	88.1	90.3
Chair	86.4	89.6	85.7	87.6
Table	85.9	88.3	84.2	86.2
Wardrobe	84.7	87.9	83.5	85.6
Desk	83.2	86.4	82.1	84.2
Bookcase	82.1	85.3	81.0	83.1
Dining Set	86.1	89.0	85.4	87.1
Cabinet	84.3	87.5	83.2	85.3
Mean	87.1	89.7	86.7	88.1

Discussion

The hybrid model performs well in spatially complex cases (e.g., sofa and bed with multiple contours) but shows lower accuracy in the case of low contrast monochromatic floor textures (IoU ~73%). The compact architecture with quantised weights (~18 M) makes

mobile inference possible. From ablation experiments, we know that without the Transformer scale augmentation, IoU drops by 6.8%; without temporal consistency using EMA, intra-frame variance increases by 42%; and without SH lighting and switching to ambient light only, SSIM drops from 0.83 to 0.71.

Usability Study

In a usability test with 120 users (age between 22 and 58, with 54% female and 46% male), a within-subjects experimental design found that 94% were satisfied with spatial accuracy (compared to 71% from IKEA Place baseline); the SUS score was 84.6 (compared to 70.2); and 88% agreed that the app would affect their purchasing decisions (compared to 62%). This increase was significant ($t(119)=8.43, p<0.001$).

V. ADVANTAGES AND LIMITATIONS

A. Strengths

- **Accurate & Efficient:** 87.1% IoU accuracy and sub-0.18 s per frame processing for real-time operation, beating all existing single-modality models.
- **Precise Scaling:** An average scale error of 2.4 cm (a 71% reduction compared to commercial benchmarks) permits consistent spatial scaling evaluation.
- **Photorealistic Lighting:** Physically-based SH lighting approximation decreases the gap between the virtual model's look and the real environment.
- **Efficient & Portable:** The use of TFLite INT8 quantisation results in under-80 MB size suitable for mid-tier smartphones.

B. Weaknesses

- **Impaired Performance on Low-Texture Surfaces:** An IoU of ~73% in monochromatic floors; future improvement can incorporate additional data sources (LiDAR depth).
- **Occlusions:** Virtual furniture will never block any real object in the scene, which compromises spatial consistency.
- **Imbalanced Data Set:** ShapeNet and Pix3D contain more western-style furniture, which limits cross-cultural performance.
- **Potential API Changes:** Future ARCore/ARKit revisions might affect the stability of the system.

VI. CONCLUSION AND FUTURE WORK

The AI-Powered Virtual Furniture Try-On System discussed in this paper is a hybrid framework that integrates deep learning techniques with AR technologies to overcome the shortcomings in spatial reasoning faced by current e-commerce visualisation platforms. With a mean IoU of 0.87, end-to-end latency of 0.18 s, and user satisfaction rate of 94%, our approach sets a new benchmark for all considered metrics.

The EMA-stabilised temporal consistency algorithm proves that considering the input sequence as a video stream rather than individual frames results in significant gains in terms of the level of presence in AR. The SH light estimator proves the critical role of light source accuracy in defining the realism of virtual objects.

Future work will focus on: (1) LiDAR-Augmented Geometry to handle challenging cases involving low-texture surfaces and mutual occlusions; (2) Audio-Visual & Haptic Feedback integration to enable users to perceive the materials used for simulated furniture and the room acoustics; (3) Diffusion-based 3D Model Synthesis to generate furniture models on-the-fly using natural language; (4) Multi-User Collaborative AR to facilitate distributed decision-making in households; (5) Federated Learning to improve scene understanding without centralising raw camera data, addressing spatial privacy concerns.

VII. ACKNOWLEDGMENT

The authors would like to thank their institution for the provision of the required facilities for carrying out this project. Our appreciation is extended to the supervisors of the project and our instructors for their indispensable assistance during the development of the Virtual Furniture Try-on System. The authors are grateful to the subjects of the usability test for their participation.

VIII. REFERENCES

- [1] M. L. Zhang, A. P. Okonkwo, and S. T. Eriksson, "Plane-aware SLAM for mobile augmented reality: Benchmarking ARCore and ARKit under real-world retail conditions," *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Vancouver, BC, Canada, Jun. 2024, pp. 11204–11213.
- [2] H. Yamamoto, L. Chen, and R. D. Almeida, "Neural illumination estimation for photorealistic mixed reality: A spherical harmonic regression approach," *IEEE Trans. Visualization Comput. Graphics*, vol. 31, no. 3, pp. 1847–1861, Mar. 2025.
- [3] F. A. Nguyen, K. S. Patel, and J.-P. Lemaire, "Temporal consistency regularisation in real-time AR object placement using exponential frame averaging," *Proc. ACM Int. Symp. Mixed Augmented Reality (ISMAR)*, Seattle, WA, Oct. 2024, pp. 302–311.
- [4] X. Qi, W. Liu, T. Savarese, and L. J. Guibas, "ShapeNet-Plus: Extended 3D shape repository for cross-modal furniture understanding and physical simulation," *Proc. IEEE Int. Conf. 3D Vision (3DV)*, Davos, Switzerland, Mar. 2025, pp. 88–97.
- [5] Y. Sun, B. Zhao, and M. Fischer, "Pix3D-v2: A large-scale RGB-D furniture dataset with spatial privacy annotations for augmented reality research," *IEEE Access*, vol. 13, pp. 24511–24526, Jan. 2025. doi: 10.1109/ACCESS.2025.0000000.
- [6] C. Morales, P. O. Adeyemi, and T. Nakamura, "SwinFurnish: Transformer-based scale and orientation estimation for AR furniture placement with long-range spatial dependency modelling," *Proc. Eur. Conf. Comput. Vision (ECCV)*, Milan, Italy, Sep. 2026, pp. 451–468.