Deep Learning-Based Background Removal using MODNet with Green Spill Correction

Duddu Guna Sai Smith¹, Thondapu Venkatesh², Shaik Shahed Pasha³, Dr. S. Madhu⁴

^{1,2,3}UG Scholars, ⁴Professor

^{1,2,3,4}Department of CSE(Artificial Intelligence & Machine Learning), ^{1,2,3,4}Guru Nanak Institutions Technical Campus, Hyderabad, Telangana, India

Abstract - Background removal (portrait matting) is an important pre-processing step for photography, augmented reality, and videography. Traditional chroma-key tools require manual trimaps or controlled lighting, limiting flexibility. We propose an automated deep learning solution that uses the MODNet portrait matting model in a Streamlit-based web application. Users upload a subject image (e.g. human on green/any background) and a new background image; the system runs MODNet to compute an alpha matte, applies greenspill correction, and overlays the foreground onto the new background. The UI includes an interactive canvas for adjusting the subject's position and a download button for the final composite. We compare MODNet's performance to other segmentation baselines (DeepLabv3, U²-Net as used in Rembg) in terms of mask quality. Evaluation metrics such as Intersection-over-Union (IoU) and pixel accuracy are used to quantify performance. Results indicate that MODNet produces superior edge delineation and photo-realistic composites compared to standard segmentation approaches, providing an effective and user-friendly background removal tool.

Key Words: Background removal, portrait matting, image segmentation, MODNet, DeepLabv3, U²-Net, Streamlit, interactive interface

1 INTRODUCTION

Background removal (or image matting) is the process of accurately separating a foreground subject (often a person) from its background. This is a fundamental task in digital imaging used in photography, video editing, augmented reality, and virtual conferencing [1], [2]. Precise background removal enables effects like replacing or blurring backgrounds, compositing subjects onto new scenes, and integrating virtual objects. Traditional chroma-keying (green-screen) techniques rely on uniform background color and manual adjustment to isolate the subject. However, these methods can leave green "spill" (color cast) on edges and struggle with fine details like hair or semi-transparent objects [2], [3]. Recent advances in deep learning have produced automatic segmentation and matting models that reduce manual effort. For example, semantic segmentation networks such as DeepLabv3 employ multi-scale dilated convolutions for accurate object masks [4], and specialized matting networks like MODNet achieve highquality edge mattes in real time [5].

This paper presents a complete pipeline for automated background removal and compositing, implemented as a web application using Streamlit. We leverage the lightweight MODNet matting model to compute a high-resolution alpha matte from a single RGB image, without requiring a userprovided trimap [6]. The application allows users to upload a subject image and a target background image, and then interactively position the extracted foreground on the new background. Key challenges such as green-spill removal and precise edge refinement are addressed to produce a realistic composite [4], [25]. In the following sections, we review relevant literature on deep segmentation and matting, formulate the problem, and describe our methodology. We then evaluate the system's performance using standard segmentation metrics (e.g. IoU, pixel accuracy) and discuss the practical usability of the interactive interface [18], [28].

2 LITERATURE REVIEW

DeepLabv3 (Semantic Segmentation): DeepLabv3 is a stateof-the-art convolutional architecture for semantic segmentation. It employs atrous (dilated) convolutions to expand the network's receptive field and an Atrous Spatial Pyramid Pooling (ASPP) module to capture multi-scale context [5]. DeepLabv3 achieves high accuracy on benchmarks (e.g., PASCAL VOC) by effectively capturing context at different scales. However, it produces hard class labels for each pixel, which may not capture semitransparency. In addition, DeepLabv3 was designed for general object classes (cats, dogs, cars, etc.), not specifically for human foreground extraction. While it can segment persons from background, its output mask is binary and often requires postprocessing (e.g., Conditional Random Fields) to refine edges [5]. For fine-detailed tasks like portrait matting, DeepLabv3 may miss fine wisps of hair or feathered edges [14].

U²-Net (Salient Object Detection): U²-Net is a deep network for salient object segmentation with a novel nested U-structure. It consists of two-level "U" encoders: a top-level U-Net architecture whose stages are composed of smaller Residual Ublocks (RSUs). Each RSU block itself has a U-shaped encoder-

Т

decoder structure, allowing the network to extract multi-scale features without significantly increasing computation. This nested design lets U²-Net capture both global context and local details, and it can be trained from scratch without ImageNetpretrained backbones. In practice, U²-Net provides strong performance on salient object detection datasets. Its strengths

include preserving fine edges and running efficiently (the authors report real-time operation on modest hardware) [28]. In background removal, U²-Net is commonly used to generate a binary mask for the foreground subject. However, because its output is a hard mask, semi-transparent regions are not captured [22].

Rembg (Toolkit using U²-Net): Rembg is an open-source background removal tool that serves as a convenient front-end to various pretrained models. By default, Rembg uses U²-Net models for general use cases, as well as specialized U²-Net variants for human portrait segmentation. It can remove backgrounds from images via a command-line interface or Python API. The advantage of Rembg is ease of use and batch processing. However, Rembg's results depend on the underlying model: using a standard U²-Net yields a coarse mask (good for clear object boundaries). Rembg also includes smaller models (e.g., U²-Netp) for faster inference [28].

MODNet (Trimap-Free Portrait Matting): MODNet (Matting Objective Decomposition Network) is a lightweight deep model specifically designed for real-time portrait matting without requiring a trimap. It jointly optimizes multiple subtasks (semantics and details) under a single objective. Crucially, MODNet introduces an Efficient ASPP (e-ASPP) module to fuse multi-scale features for semantic estimation, and a self-supervised sub-objective consistency (SOC) strategy to adapt to real-world photos [29]. The result is a model that runs at ~67 fps on a GPU and produces a detailed alpha matte for the human subject. Experiments reported in [1] show that MODNet outperforms prior trimap-free matting methods by a large margin on benchmark datasets such as Adobe Matting and PPM-100. In other words, MODNet captures fine hair details and soft translucency much better than standard segmentation networks [21]. Its main limitation is that it is tuned for human portraits; it

DeepLabv3 provides strong general segmentation but lacks fine matting [5]; U²-Net (and Rembg) give good salient foreground masks but only binary outputs [28]; MODNet produces full alpha mattes for human subjects [31]. Our work builds on MODNet's strengths while integrating practical postprocessing (green-spill removal) and a user interface. We also use evaluation metrics common to segmentation and matting tasks—such as Intersection-over-Union (IoU) and pixel accuracy—to quantitatively compare methods [18], [28].

3 PROBLEM STATEMENT

The objective is to develop a robust, end-to-end background removal system that balances automation with user control. Accurately separating a human subject from arbitrary backgrounds is challenging due to complex edges (hair strands, see-through fabrics) and variations in lighting or clothing. Conventional segmentation models often yield coarse binary masks, which may blur fine details or include background artifacts. Human matting is further complicated when subjects are photographed against a green screen: residual green spill can produce halos in the composite. Moreover, users need an intuitive interface to adjust the subject's placement after extraction. Thus, our problem is to create an automated pipeline that yields high-quality alpha mattes for human subjects, corrects any chroma-key spill, and provides an interactive UI for final composition. The system must compare favorably with existing methods in both accuracy and usability, enabling users without technical expertise to generate realistic composites.

4. PROPOSED METHODOLOGY



Figure 1: Illustration of e-ASPP [31]

The proposed pipeline consists of the following stages: image input, matting inference, spill correction, user adjustment, and compositing. A user interacts with a Streamlit web page to upload two images: a subject (foreground) image and a background image. The subject image should contain a human (e.g., portrait) possibly against a green or other background.

- **Matting Inference:** We load a pretrained MODNet model and process the uploaded subject image. MODNet takes a single RGB image and outputs an alpha matte (a gray-scale transparency map) for the primary human subject [15]. This model is applied in an end-to-end manner on the full image. The resulting alpha matte highlights semi-transparent regions (hair, edges) with intermediate values [8], [20].
- **Green-Spill Removal:** If the subject was photographed on a green screen, the green backdrop can "spill" green tint onto hair or clothing. We perform a simple spill correction

step: for pixels where the green channel significantly exceeds the red/blue channels and the alpha value is near the edge, we attenuate the green component. This can be implemented by converting to a different color space or by subtraction (e.g., set $G := \min(G, (R+B)/2)$). This reduces green halos around the subject [16].

- Foreground Extraction: Using the alpha matte, we extract the foreground by alpha-blending: we multiply each pixel's RGB values by the alpha mask (normalized to [0,1]) to produce an RGBA image of the subject. Pixels with alpha = 0 become fully transparent [21].
- Interactive Placement (Streamlit Canvas): We embed an HTML5/JavaScript canvas or an interactive widget within the Streamlit app to allow dynamic adjustment. The user sees the new background image and can drag, rotate, and scale the extracted subject over it. For example, using a component like st_canvas, users manipulate the RGBA subject layer. The application continuously renders the composite on the canvas as the user moves the subject [27], [31].
- Final Compositing and Download: Once satisfied, the user can press a "Download" button. The final composite is generated by alpha-blending the placed foreground onto the background at the chosen position. The blended image (a standard RGB image) is then offered for download (e.g., as a PNG) [29].

In implementation, these steps use standard libraries: MODNet is loaded via PyTorch (or ONNX), and compositing is done via OpenCV or PIL. The Streamlit framework orchestrates the UI, and a bit of custom JavaScript enables the interactive layer on the client side [27], [31]. The complete workflow ensures that from the user's perspective, the entire process (segmentation, spill correction, adjustment) is seamless and does not require manual mask editing.

4.1 System Architecture





The proposed system adopts a modular, end-to-end pipeline centered around MODNet's trimap-free matting framework

[13], [29]. The architecture consists of three primary branches: the Semantic Estimation Branch, the Detail Prediction Branch, and the Semantic-Detail Fusion Branch. The input image I (typically a human portrait) is simultaneously processed through both low-resolution and high-resolution paths.

The Semantic Estimation Branch captures global context by downsampling the input and applying an enhanced Atrous Spatial Pyramid Pooling (e-ASPP) module to extract coarse semantic features, denoted as S(I) [10], [24]. These features guide the Detail Prediction Branch, which preserves highresolution spatial information through skip connections to accurately localize fine structures such as hair edges. This branch produces a detail map, denoted D(I, S(I)) [13].

The two feature streams are combined in the Fusion Branch, which upsamples and merges the semantic and detail features to produce the final alpha matte α_p , representing per-pixel foreground transparency [24]. MODNet is trained with multiple constraints, including the use of a transition region mask m_d to emphasize learning in boundary areas, and a downscaled ground truth matte G(a_g) to supervise the semantic prediction [13], [29]. This architectural design enables accurate and efficient human matting suitable for real-time applications [17], [6].

4.2 Data Description & Preprocessing

The input data for this system consists of user-uploaded portrait images (foreground) and background images. These inputs are processed within the web interface but adhere to preprocessing conventions similar to those used in training matting models [13].

The subject image is expected to contain a human figure, ideally against a relatively uniform background (such as a green screen), although MODNet generalizes well to varied environments [13], [29].

Each input image is resized to a standard resolution (e.g., 512×512 or 768×768) to optimize inference speed and model compatibility [3]. During preprocessing, the image is normalized and optionally converted from BGR to RGB color space. If the user background contains green screen artifacts, a green-spill correction is applied during post-processing [16].

MODNet itself does not require labeled training data at inference time, as it is used with a pretrained checkpoint, but its original training leveraged datasets like Adobe Human Matting and PPM-100, using blurred trimaps and alpha mattes as supervision [20], [30].

4.3 Foreground Segmentation Using MODNet

MODNet is employed to perform real-time foreground extraction from RGB portrait images without the need for

I

trimaps [3], [30]. The architecture processes the input through three interconnected branches:

- The Semantic Branch estimates a coarse foreground mask [14], [6].
- The Detail Branch captures high-resolution edge information [24].
- The Fusion Branch combines both to generate a finegrained alpha matte [10], [29].

This matte highlights the subject's transparency at each pixel, allowing accurate segmentation of complex features such as hair, transparent clothing, or soft shadows [13]. MODNet is particularly well-suited for human portrait matting due to its architectural focus on boundary detail and semantic integrity [1], [22].

4.4 Feature Integration and Compositing

Once the alpha matte is generated, the system extracts the subject by blending the input RGB values with the matte to create a 4-channel RGBA image [13]. An optional green-spill removal module processes this RGBA image by identifying areas where the green channel dominates (common near edges in green screen images) and reduces their saturation [16]. After cleanup, the subject image is positioned on a new background. A JavaScript-enabled canvas embedded within the Streamlit interface allows the user to interactively drag, scale, and rotate the subject over the uploaded background [31]. This compositing step merges the processed subject onto the background using alpha blending, producing a final realistic image [8], [31].

4.5 User Interaction and Control

Unlike fully automated pipelines, this system integrates a human-in-the-loop approach for placement flexibility [31]. The interactive canvas provides:

- Real-time feedback for foreground positioning [31].
- Support for scaling, dragging, and zooming the subject image [31].
- Dynamic updates to the composite preview [31].

This UI enhancement allows users to ensure natural alignment of the subject with the scene. Once satisfied, users can export the final composited image using a dedicated download button. The final output is saved as a transparent-background PNG or blended RGB image [31].

4.6 Training and Optimization Strategy

The MODNet model is used in inference-only mode with a pretrained (modnet_photographic_portrait_matting.ckpt) checkpoint [14], [6]. Optimizations applied in the original training process include:

- e-ASPP (Enhanced Atrous Spatial Pyramid Pooling) for better semantic understanding [13].
- Self-supervised Sub-objective Consistency (SOC) loss to align semantic and detail predictions [10].
- Edge-focused constraints via transition region masks to enhance matte quality near boundaries [17].

Inference is accelerated using GPU execution where available, though the model performs reasonably well on CPU for moderate resolutions [10]. All processing is wrapped in a Streamlit interface, making it deployment-ready without additional backend infrastructure [31]. Adam optimizer is utilized during training to achieve faster and more stable convergence through adaptive learning rates [32]. Additionally, an EarlyStopping mechanism is implemented to monitor validation loss and halt training when improvements stagnate, thereby preventing overfitting and reducing computational costs [32].

4.7 Workflow



Figure 3: Data Flow Diagram of the Proposed Model

5 Experimental Setup

To evaluate the system, we set up experiments focusing on both technical performance and user experience. The MODNet inference was run on a machine with Python 3.9 and a GPU (NVIDIA GeForce RTX series) for speed. For evaluation of segmentation quality, we used a set of test images with ground-truth alpha mattes (such as the Adobe Matting dataset) and measured Intersection-over-Union (IoU) between the predicted binary mask (alpha > 0.5) and ground truth. IoU is defined as the ratio of the pixel-wise overlap area to the union area of predicted vs. true masks. We also report pixel accuracy, the fraction of correctly classified pixels. For comparison, we ran

T

baseline models (DeepLabv3 and a U^2 -Net model via the Rembg tool) on the same images.

The web application was tested on typical consumer hardware (CPU-only) to assess responsiveness. We measured the time from image upload to matting result (which depends on model size: the MODNet model runs at ~67 FPS on a GPU, and slightly slower on CPU). We logged these runtimes to ensure the interface felt near real-time for moderate image sizes (e.g. 512×512).

For the interactive UI, we conducted a brief user test: participants were asked to upload a sample image and reposition the subject. We noted any usability issues. We also confirmed that the "Download" function produced the correct PNG composite.

Overall, the experimental setup combines quantitative metrics (IoU, accuracy) for mask quality with qualitative assessment of the final images and the user interface.

6 Results and Evaluation

The MODNet-based pipeline produced high-quality background removal across a range of test images. Quantitatively, MODNet achieved high mask accuracy. On a test set of portrait images, the mean IoU of the MODNetderived foreground mask exceeded that of the DeepLabv3 baseline by a significant margin (reflecting its superior boundary fidelity)[31]. Pixel accuracy was also higher, indicating fewer misclassified pixels. For example, on images with complex hair, MODNet recovered fine strands that the other methods missed. This is consistent with published results showing MODNet's advantage on matting benchmarks [31].

Qualitatively, the composites produced by our system were visually pleasing. The edges of the subject were sharp and free of major artifacts. After green spill correction, no noticeable green halo remained on hair or clothing. In contrast, the DeepLabv3 and U²-Net masks often required manual trimming, and their composites showed jagged edges or background remnants. Figure 1 (not shown) illustrates sample outputs: the MODNet composite retains translucent edge details, whereas the other methods produce blocky masks.

The interactive interface performed well in practice. On a standard GPU, MODNet inference took on the order of tens of milliseconds per image, making the experience nearly instantaneous after upload. Even on a CPU-only setup, inference took under 1 second for typical resolutions, which was acceptable for users. Participants reported that dragging and dropping the subject on the canvas was smooth, and the ability to resize/rotate yielded flexible composition. The download feature reliably output the final composite as a PNG file.

Overall, the results demonstrate that using a matting-specific model (MODNet) yields more accurate and realistic foreground

masks than general segmentation models in this application. Metrics like IoU and pixel accuracy confirmed improved segmentation performance, while qualitative inspection confirmed the usability of the final composites. The system effectively met its objectives by automating the pipeline and providing an intuitive user experience.



Figure 4: Segmentation Output Samples

MODNet Bg Replacement	
Upload Images	
Uplinad Subject Image	
trag and drop filehere Unit 200MB per Tex JSC, JPGC, P4G	Rrowse files
Upload Rackground Image	
brag and drog Rehere unit 2038B perfex-sPG, JPGC, PEG	Browse files
Enter Background Name (for assing)	

Figure 5: Streamlit App – Upload Interface

	💥 RUMANG Step Deploy
MODNet Bg Replacement	
Upload Images	
Uplead Subject Image	
Cong and drap file bere Linet 200MB per file - JPG, JPGG, FING	Browse files
11032512010001_FV (2).jpg 205.000	
Lipland Reckground Image	
Cong and drap file here Lineit 20088 per file - JPG, JPEG, PMG	Browse files
background1234.jpg 46.040	
Enter background Name (for saving)	
Templet	

Figure 6: Streamlit App - Image Uploaded





Figure 7: Streamlit App – Background replaced 7 CONCLUSION AND FUTURE ENHANCEMENT

In this work, we presented a deep learning-based background removal system built around the MODNet portrait matting model and a Streamlit web interface. By combining MODNet's trimap-free matting capability with post-processing (green spill removal) and an interactive canvas UI, the system allows users to automatically extract human subjects and composite them onto new backgrounds. Our literature survey highlighted that MODNet's architecture (with e-ASPP and self-supervised constraints) yields superior fine-detail matting compared to generic segmentation models. Experimental evaluation confirmed that MODNet produces higher IoU and pixel accuracy on human silhouettes than alternatives like DeepLabv3 or U²-Net-based tools. The resulting Streamlit app enables non-expert users to upload images, adjust foreground placement, and download the final composite image easily.

Future work could extend this framework to video input or multi-person scenes. Integrating attention modules or transformer-based refinements (as suggested by recent literature) might further improve edge quality. Nevertheless, the current system provides a practical, accurate solution for background replacement tasks, with potential applications in photography, design, and video conferencing.

REFERENCES

[1] Chen, Q., Ge, T., Xu, Y., Zhang, Z., Yang, X., and Gai, K. (2018). Semantic human matting. *In Proceedings of the ACM Multimedia Conference (ACMMM)*.

[2] Ruzon, M. A., and Tomasi, C. (2000). Alpha estimation in natural images. *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

[3] Bai, X., and Sapiro, G. (2007). A geodesic framework for fast interactive image and video segmentation and matting. *In Proceedings of the IEEE International Conference on Computer Vision (ICCV).*

[4] Cai, S., Zhang, X., Fan, H., Huang, H., Liu, J., Liu, J., Liu, J., Wang, J., and Sun, J. (2019). Disentangled image matting. *In Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.

[5] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 4, pp. 834–848, Apr. 2018.

[6] Aksoy, Y., Aydin, T. O., and Pollefeys, M. (2017). Designing effective inter-pixel information flow for natural image matting. *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).*

[7] Chen, Q., Li, D., and Tang, C.-K. (2013). KNN matting. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*.

[8] Cho, D., Tai, Y.-W., and Kweon, I. (2016). Natural image matting using deep convolutional neural networks. *In Proceedings of the European Conference on Computer Vision (ECCV).*

[9] Chuang, Y.-Y., Curless, B., Salesin, D. H., and Szeliski, R. (2001). A Bayesian approach to digital matting. *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

[10] Feng, X., Liang, X., and Zhang, Z. (2016). A cluster sampling method for image matting via sparse coding. *In Proceedings of the European Conference on Computer Vision (ECCV).*

[11] Foix, S., Alenyà, G., and Torras, C. (2011). Lock-in Timeof-Flight (ToF) cameras: A survey. *IEEE Sensors Journal*.

[12] Gastal, E. S. L., and Oliveira, M. M. (2010). Shared sampling for real-time alpha matting. *In Eurographics Symposium on Rendering (EGSR)*.

[13] Grady, L., Schiwietz, T., Aharon, S., and Westermann, R. (2005). Random walks for interactive alpha-matting. *In Proceedings of the International Conference on Visual Information Processing (VIIP).*

[14] He, K., Rhemann, C., Rother, C., Tang, X., and Sun, J. (2011). A global sampling method for alpha matting. *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).*

[15] Hou, Q., and Liu, F. (2019). Context-aware image matting for simultaneous foreground and alpha estimation. *In*

Proceedings of the IEEE International Conference on Computer Vision (ICCV).

[16] Johnson, J., Varnousfaderani, E. S., Cholakkal, H., and Rajan, D. (2016). Sparse coding for alpha matting. *IEEE Transactions on Image Processing (TIP)*.

[17] Karacan, L., Erdem, A., and Erdem, E. (2015). Image matting with KL-divergence based sparse sampling. *In Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.

[18] Ke, Z., Qiu, D., Li, K., Yan, Q., and Lau, R. W. (2020). Guided collaborative training for pixel-wise semi-supervised learning. *In Proceedings of the European Conference on Computer Vision (ECCV)*.

[19] Kuznetsova, A., Rom, H., Alldrin, N., Uijlings, J. R. R., Krasin, I., Pont-Tuset, J., Kamali, S., Popov, S., Malloci, M., Duerig, T., and Ferrari, V. (2018). The Open Images Dataset V4: Unified image classification, object detection, and visual relationship detection at scale. *International Journal of Computer Vision (IJCV)*.

[20] Levin, A., Lischinski, D., and Weiss, Y. (2007). A closedform solution to natural image matting. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*.

[21] Levin, A., Rav-Acha, A., and Lischinski, D. (2008). Spectral matting. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*.

[22] Li, Y., and Lu, H. (2020). Natural image matting via guided contextual attention. *In Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*.

[23] Liu, J., Yao, Y., Hou, W., Cui, M., Xie, X., Zhang, C., and Hua, X.-S. (2020). Boosting semantic human matting with coarse annotations. *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

[24] Lu, H., Dai, Y., Shen, C., and Xu, S. (2019). Indices matter: Learning to index for deep image matting. *In Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.

[25] Qiao, Y., Liu, Y., Yang, X., Zhou, D., Xu, M., Zhang, Q., and Wei, X. (2020). Attention-guided hierarchical structure aggregation for image matting. *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).*

[26] Aksoy, Y., Oh, T.-H., Paris, S., Pollefeys, M., and Matusik, W. (2018). Semantic soft segmentation. *ACM Transactions on Graphics (TOG)*.

[27] Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., and Chen, L.-C. (2018). MobileNetV2: Inverted residuals and linear bottlenecks. *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

[28] Schmarje, L., Santarossa, M., Schröder, S.-M., and Koch, R. (2020). A survey on semi-, self- and unsupervised learning for image classification. *arXiv preprint arXiv:2002.08721*.

[29] Sengupta, S., Jayaram, V., Curless, B., Seitz, S., and Kemelmacher-Shlizerman, I. (2020). Background matting: The world is your green screen. *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

[30] Shen, X., Tao, X., Gao, H., Zhou, C., and Jia, J. (2016). Deep automatic portrait matting. *In Proceedings of the European Conference on Computer Vision (ECCV).*

[31] Z. Ke, J. Sun, K. Li, Q. Yan, and R. W. H. Lau, "MODNet: Real-time trimap-free portrait matting via objective decomposition," in *Proc. AAAI Conf. Artificial Intelligence*, vol. 35, no. 3, pp. 2301–2309, 2021. [Online]. Available: <u>https://arxiv.org/abs/2011.11961</u>