

Nature Unveiled: Recovering Occluded Landscape from Video Sequence

Anonymous CVPR submission

Paper ID *****

Abstract

Background recovery in videos of natural scenes is highly desirable given the increased movement and clutter of human activity in natural spaces. This study aims to remove occlusion and recover the natural backdrop from a video sequence. The process begins by segmenting the video frames into foreground and background. Among various segmentation methods tested, using learned optical flow between frames yielded the best results. After calculating a segmentation/mask for each frame, the next step is to inpaint the foreground. Inpainting methods tested include Stable Diffusion Inpainting and Flow-edge Guided Video Completion (FGVC). While Stable Diffusion Inpainting produced the highest quality frames, it ultimately lacked consistency between frames. The discussion includes potential future work to ensure frame-to-frame consistency.

1. Introduction

Background recovery in videos of natural scenes is highly desirable due to the increased movement and clutter of human activity in these spaces. Objects such as telephone lines, fences, houses, and trees can obscure views of scenic landscapes. This project aims to recover the background and remove occlusions from a video sequence shot from the passenger side of a car, leveraging information from neighboring frames to "see" behind occlusions.

1.1. Previous Work

Previous work on background recovery from video includes several approaches. For instance, [8] focuses on using the estimated motion of moving occlusions to recover background scenes. However, this method assumes a stationary camera and relies on multiple frames where the occluded object is moving to reconstruct the background behind it.

In [9], the authors leverage motion parallax from slight camera movements and accompanying optical flow fields to separate foreground and background. This approach successfully removes reflections and fence occlusions through

a unified framework, demonstrating the potential of combining slight camera motion with advanced optical flow techniques.

Another significant contribution is from [7], where the authors aim to remove occluded objects such as cars, pillars, and people to clean up photos for Image-Based Rendering (IBR). They use multiple images to create a rough 3D rendering of a scene and then employ this 3D model to inpaint masked regions, thereby enhancing the visual quality and coherence of the reconstructed scenes.

This work builds on these approaches by focusing on dynamic scenes with a moving camera, aiming to improve frame consistency and the quality of the recovered background. Future work will explore methods to ensure temporal consistency between frames, addressing the challenges identified with current inpainting techniques.

Here is a revised and expanded version of the Problem Statement section:

2. Problem Statement

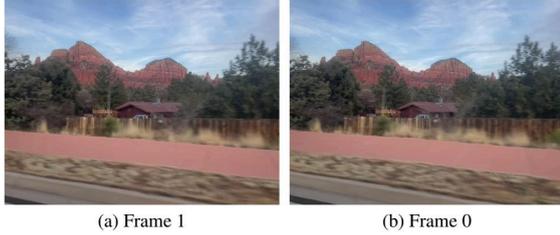
Assume that there is a background B that is occluded by a foreground F . We are given multiple continuous frames in a video where we can estimate the depth to separate the objects in B and F . Due to the changing position of the camera and the non-uniform and semi-transparent nature of F in some areas, we can accumulate information about B over time.

The goal is to use this accumulated information to inpaint the photos in a manner that is consistent across views. We aim to compare this multi-frame inpainting method with single-frame inpainting results to evaluate improvements in consistency and quality.

2.1. Dataset

Finding an appropriate synthetic dataset with ground truth data for this task proved challenging. Most available landscape or driving footage, both synthetic and real, is typically captured with the camera mounted on the front of the vehicle rather than to the side, which is necessary for our specific occlusion recovery problem.

108 Currently, I am testing my method on a video sequence
109 that I captured. This sequence is 4 seconds long and consists
110 of 120 frames moving from right to left. Frames 0 and 1 are
111 shown in Figure 1.
112



122 Figure 1. Video frames with occluded landscape. Notice the camera
123 is moving to the left.
124

125 The dataset for this project is a custom video sequence
126 shot from the passenger side of a car. This allows for dynamic
127 occlusion scenarios where the camera moves and captures varying
128 perspectives of the occluded background. The sequence captures the
129 complexity of real-world occlusions and provides a challenging
130 testbed for evaluating the proposed background recovery method.
131
132

133 3. Technical Approach

134 To solve the problem of background recovery in videos
135 with occlusions and moving cameras, the following technical
136 approach will be employed:
137

138 3.1. Segmentation

139 The first step is to segment the frames of the video into
140 foreground (F) and background (B). Accurate segmentation
141 is crucial for isolating the regions that need inpainting. The
142 segmentation process relies on calculating a depth map for
143 each frame. This can be done using:
144

- 145 • **Stereo Depth Map (neighboring frames):** Utilizing
146 the disparity between consecutive frames to estimate
147 object depth.
148
- 149 • **Monocular Depth Estimation (single frame):** Using
150 deep learning models trained on large datasets to predict
151 depth from a single frame.
152
- 153 • **Learned Optical Flow:** Leveraging parallax effect
154 (objects far away move less across frames) and using
155 optical flow between frames as an indicator of object
156 depth.

157 3.2. Masked Inpainting

158 Once the foreground regions are identified and masked,
159 inpainting techniques are applied to recover the back-
160 ground:
161

- 162 • **Text-guided Inpainting with Stable Diffusion:** Utilizing
163 Stable Diffusion models to inpaint masked regions
164 guided by text prompts. This method focuses
165 on generating high-quality images but needs to address
166 temporal consistency across frames.
167
- 168 • **Flow-edge Guided Video Completion (FGVC):** Using
169 optical flow and edge information to guide the in-
170 painting process, ensuring that the inpainted regions
171 are temporally coherent and consistent with the video
172 structure.
173

174 3.3. Comparison of Methods

175 The performance of the inpainting methods will be com-
176 pared based on several criteria:
177

- 178 • **Visual Quality:** Evaluating the perceptual quality of
179 the inpainted frames for naturalness and realism.
180
- 181 • **Consistency:** Assessing the temporal consistency
182 across frames to avoid flickering and artifacts.
183

184 4. Segmentation Results

185 4.1. Stereo Depth Map Estimation

186 I experimented with stereo methods using OpenCV [1].
187 Given that the camera is moving along a straight trajectory,
188 the resulting frames are nearly parallel, which facilitates the
189 use of stereo vision techniques. The depth map generated
190 from the disparities between frame 0 and frame 1 is shown
191 in Figure 2, along with the corresponding segmented fore-
192 ground and background images in Figure 3
193

194 For the trees and fence, the depth map performs quite
195 well, accurately distinguishing these elements. However,
196 for the sidewalk, grass, and curb, the uniform texture
197 results in an incorrect depth map in those regions. This
198 highlights a common challenge in stereo depth estimation
199 where regions with uniform texture or repetitive patterns
200 can lead to ambiguous disparity calculations.
201

202 4.1.1 Find Corresponding Points

203 I wondered if the depth map from stereo would be improved
204 by rectifying the images since the images are not parallel.
205 I used OpenCV [1] to find SIFT (Scale-Invariant Feature
206 Transform) points and brute force block match 150 points
207 on the image. The corresponding points and SIFT points
208 are shown in Figure 4.
209

210 4.1.2 Rectify the Images w/ Corresponding Points

211 Using HW 2 Problem 2, I calculated the homographies to
212 create parallel images as seen in Figure 5. However, now
213 the parallel images are different dimensions and OpenCV
214 function only take in images of the same dimensions.
215

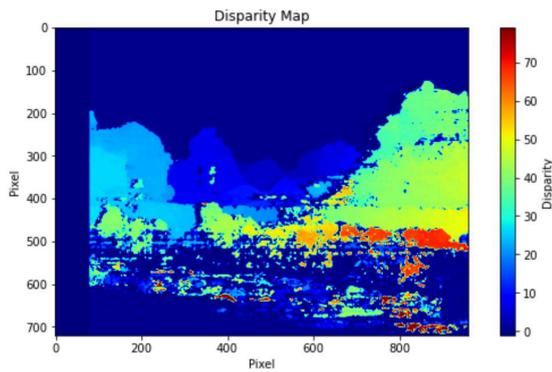


Figure 2. OpenCV Stereo calculated depth map from frame 0 and frame 1.

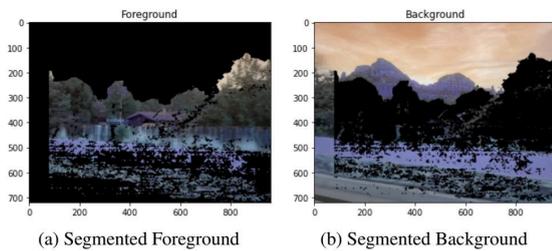


Figure 3. Foreground and background segmented with stereo depth map.

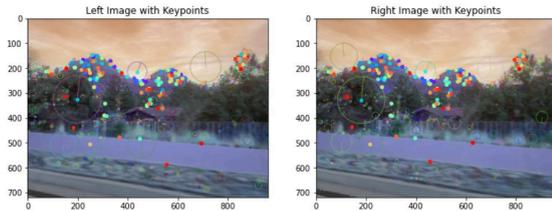


Figure 4. Corresponding points and SIFT points.

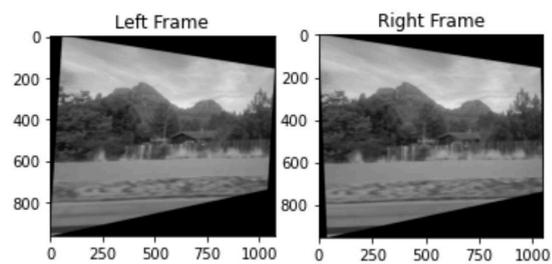


Figure 5. Rectified parallel images of frame 0 and frame 1.

4.2. Monocular Depth Estimation

Given the difficulties of stereo depth estimation, I instead turned to estimating depth from a single frame.

4.2.1 GLPN fine-tuned on NYUv2

I used a pretrained global-local path networks (GLPN) model [3] and the predicted depth map for Frame 0 is seen in Figure 6. The depth map is too smooth and notice the tree (circled in red) is blended out of the depth map. The resulting foreground and background segmentation is shown in Figure 7. Notice the house in the middle and trees behind it are not segmented correctly.

4.2.2 MiDaS

I tried another pretrained model MiDaS [4]. The estimated depth map of Frame 0 is shown in Figure 8 and corresponding foreground and background segmentation is shown in Figure 9. The model was successful in predicting the house in the middle to be on the same depth plane. However this depth map is too smooth to segment the trees in the middle into the foreground.

The subpar performance of monocular depth estimation led me to look at learned optical flow methods instead.

4.3. Learned Optical Flow

I used a pretrained optical flow model RAFT [6] in the Pytorch library. The optical flow for each frame was calculated with its following frame. Since objects that are closer to the moving camera will move at a faster speed than objects away from the camera, we can use the magnitude of the optical flow as an indicator of objects' depth in the scene. A random selection of frames with corresponding optical flow magnitude maps are shown in Fig. 10.

From the magnitude maps of each frame, I then looked at the best threshold to use for the foreground mask. An example of various threshold levels (relative to max value of optical flow magnitude of a frame) are shown in Fig. 11. The final masks that were used for inpainting were magnitude greater than 5 % of the max value of the optical flow magnitude.

5. Inpainting Results

5.1. Stable Diffusion Inpainting

Using a pretrained model from Hugging Face that is based off of Stable Diffusion that takes in a text prompt, image, corresponding mask and returns an image with the mask region inpainted according to the text prompt [5]. An example of outputs to the prompt "mountains in the desert" is shown in Fig. 12.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377



Figure 6. GLPN Monocular Depth Map on Frame 0

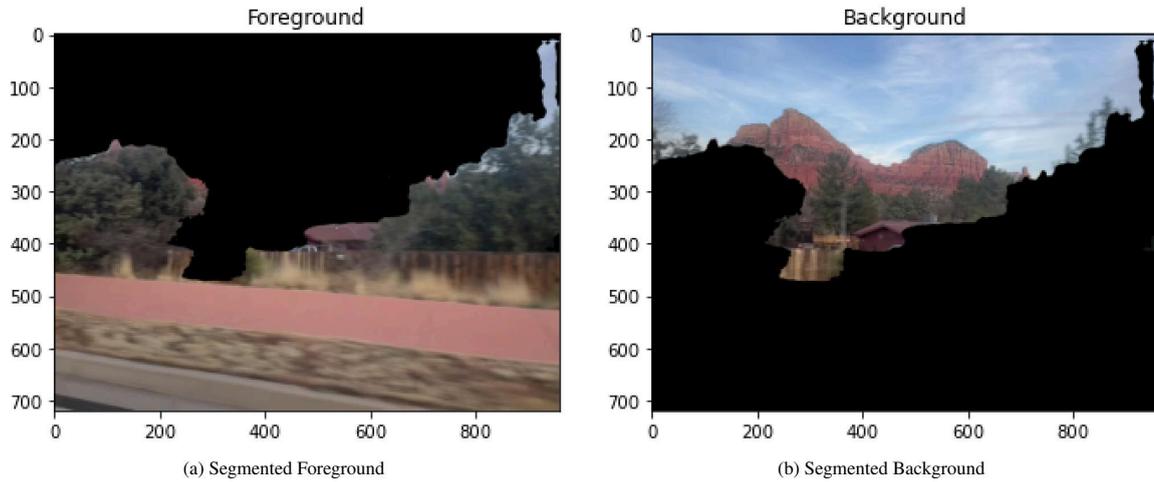


Figure 7. Foreground and background segmented with GLPN network.

Since the pretrained model processed each frame independently from another the resulting frames from a video were not consistent despite using the same text-prompts and random seeds as shown in Fig. 13.

5.2. Flow-edge Guided Video Completion

I then fed frames and corresponding masks into a flow-based video completion algorithm called "Flow-edge Guided Video Completion" (FGVC) [2]. This algorithm relies on optical flow to fill in masked regions across frames.

An resulting frame is shown in Fig. 14. This method is normally used with mask of smaller area (e.g. removing a tennis player in a tennis match) and fails to find the correct information to fill in the large mask region of the frames.

6. Conclusion

Among the methods tested, optical flow proved to be the most effective for depth level segmentation compared to monocular and stereo depth estimation. It provided the most

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

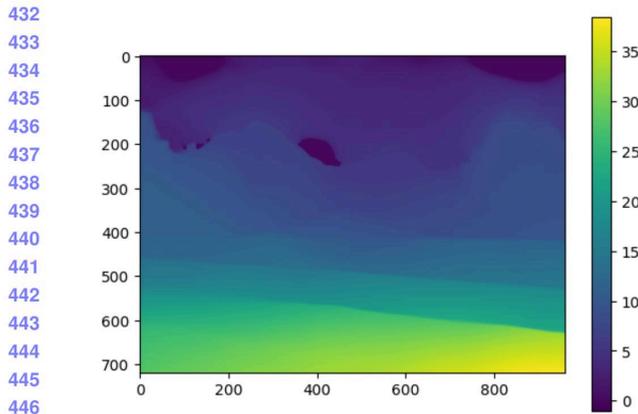


Figure 8. MiDaS Monocular Depth Map on Frame 0

accurate distinction between foreground and background objects.

However, inpainting using diffusion models, while capable of producing high-quality frames, requires specific prompts and lacks consistency across all frames even when using the same prompt. The Flow-edge Guided Video Completion (FGVC) method also fell short, as the regions to be inpainted were too large and there wasn't sufficient information across frames to accurately complete the masked regions.

To address the issue of inconsistent inpainting, a potential solution could involve training a Low-Rank Adaptation (LORA) on the video or fine-tuning the existing diffusion inpainting model. This approach would allow the model to take a prior image into account and penalize deviations from it, potentially improving temporal coherence and consistency in the inpainting results.

Future work will focus on these enhancements to improve the robustness and reliability of background recovery in video sequences.

References

- [1] G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000. 2
- [2] Chen Gao, Ayush Saraf, Jia-Bin Huang, and Johannes Kopf. Flow-edge guided video completion. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XII 16*, pages 713–729. Springer, 2020. 4
- [3] Doyeon Kim, Woonghyun Ga, Pyunghwan Ahn, Donggyu Joo, Sehwan Chun, and Junmo Kim. Global-local path networks for monocular depth estimation with vertical cutdepth. *CoRR*, abs/2201.07436, 2022. 3
- [4] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular

- depth estimation: Mixing datasets for zero-shot cross-dataset transfer, 2020. 3
- [5] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, June 2022. 3
- [6] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 402–419. Springer, 2020. 3
- [7] Theo Thonat, Eli Shechtman, Sylvain Paris, and George Drettak. Multi-view inpainting for image-based scene editing and rendering. In *2016 Fourth International Conference on 3D Vision (3DV)*, pages 351–359. IEEE, 1
- [8] Srenivas Varadarajan, Lina J. Karam, and Dinei Florencio. Background recovery from video sequences using motion parameters. In *2009 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 989–992. ISSN: 2379-190X. 1
- [9] Tianfan Xue, Michael Rubinstein, Ce Liu, and William T. Freeman. A computational approach for obstruction-free photography. 34(4):1–11. 1

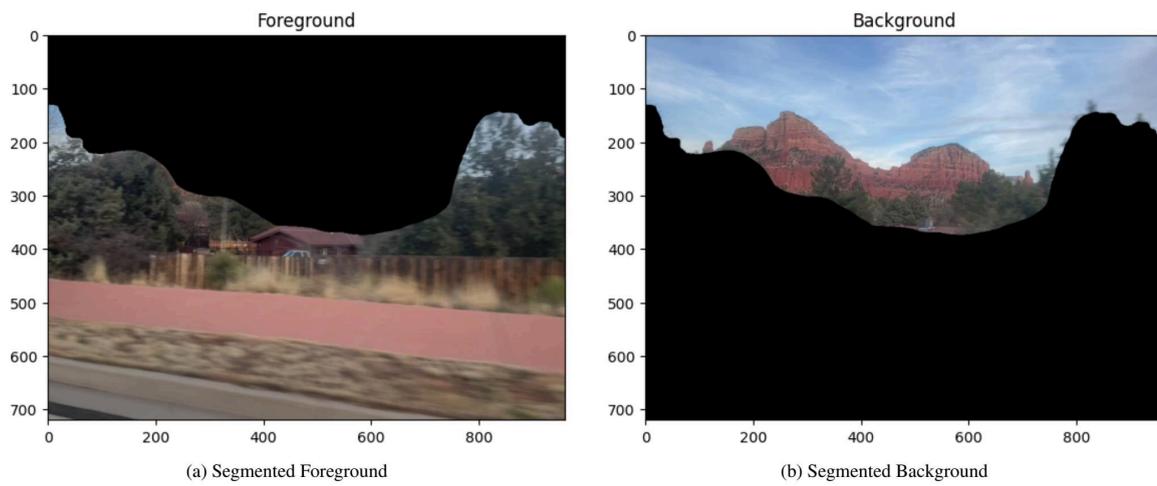


Figure 9. Foreground and background segmented with MiDaS network.

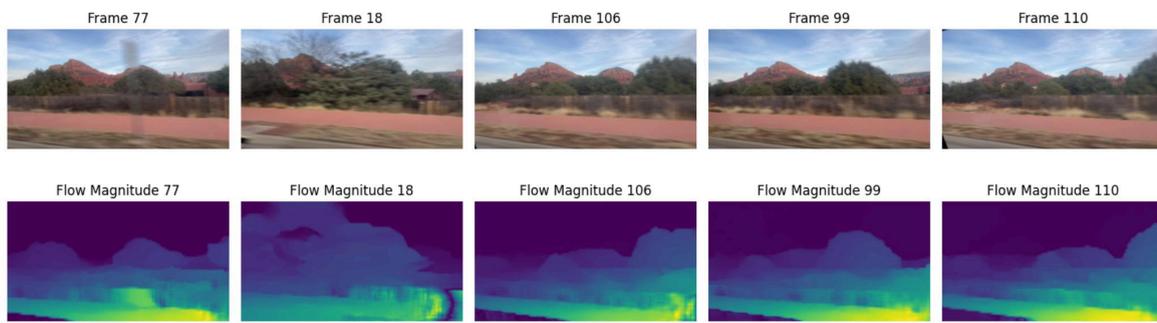


Figure 10. Optical Flow Magnitude for Random Frames from RAFT

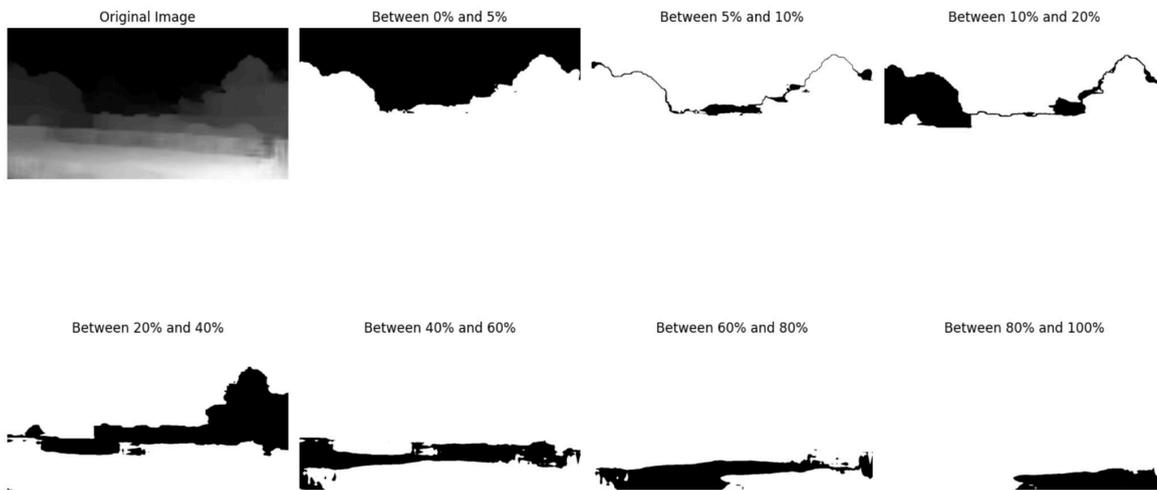


Figure 11. Example of Threshold Layers of Optical Flow Magnitude

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755



Figure 12. Diffusion model outputs of masked image with prompt "mountains in the desert"



Figure 13. Selection of Diffusion Model Inpainted Frames with "desert landscape" text prompt

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

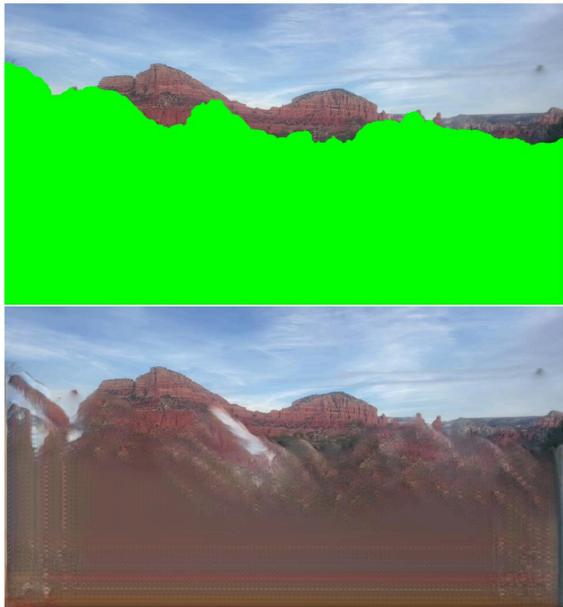


Figure 14. Image with mask applied and corresponding FGVC calculated frame

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863