

Self-Supervised Video Understanding

Erika Lu

Kellogg College
University of Oxford

*A thesis submitted for the degree of
Doctor of Philosophy*

Trinity 2021

Abstract

The advent of deep learning has brought about great progress on many fundamental computer vision tasks such as classification, detection, and segmentation, which describe the categories and locations of objects in images and video. There has also been much work done on supervised learning—teaching machines to solve these tasks using human-annotated labels. However, it is insufficient for machines to know only the names and locations of certain objects; many tasks require a deeper understanding of the complex physical world—how objects interact with their surroundings, for example (often by creating shadows, reflections, surface deformations, and other visual effects). Furthermore, training models to solve these tasks while relying heavily on human supervision is costly and impractical to scale. Thus, this thesis explores two directions: first, we aim to go beyond segmentation and address a wholly new task: grouping objects with their correlated visual effects (e.g. shadows, reflections, or attached objects); second, we address the fundamental task of video object segmentation in a self-supervised manner, without relying on any human annotation.

To automatically group objects with their correlated visual effects, we adopt a layered approach: we aim to decompose a video into object-specific layers which contain all elements moving with the object. One application of these layers is that they can be recombined in new ways to produce a highly realistic, altered version of the original video (e.g. removing or duplicating objects, or changing the timing of their motions). Here the key is to leverage natural properties of convolutional neural networks to obtain a layered decomposition of the input video. We design a neural network that outputs layers for a video by overfitting to the video. We first introduce a human-specific method, then show how it can be adapted to arbitrary object classes, such as animals or cars.

Our second task is video object segmentation: producing pixel-wise labels (segments) for objects in videos. Whereas our previous work is optimized on a single video, here we take a data-driven approach and train on a large corpus of videos in a self-supervised manner. We consider two different task settings: (1) semi-supervised object segmentation, where an initial object mask is provided for a single frame and the method must propagate this mask to the remaining frames, and (2) moving object discovery, where no mask is given and the method must segment the salient moving object. We explore two different input streams: RGB and optical flow, and discuss their connection to the human visual system.

Self-Supervised Video Understanding



Erika Lu
Kellogg College
University of Oxford

A thesis submitted for the degree of
Doctor of Philosophy

Trinity 2021

Acknowledgements

I am extremely grateful to my supervisor, Professor Andrew Zisserman, without whom this thesis would not have been possible. His support and guidance were integral to my growth during this PhD, and I am truly thankful to have had the opportunity to learn from him.

I would also like to give a heartfelt thanks to Weidi Xie for his mentorship. This PhD would not have been successful without his kindness and support. I am also extremely grateful to my mentors and collaborators at Google: Bill Freeman, Miki Rubinstein, Forrester Cole, Tali Dekel, and David Salesin. I learned a great deal from them, and it was a pleasure to work with them.

A big thank you to my collaborators in VGG: Zihang Lai, Charig Yang, and Hala Lamdouar. I would like to thank as well all of the other members of VGG (who are too many to name) for providing a friendly and supportive atmosphere, and for sharing many meals together. I gratefully acknowledge the EPSRC and Oxford-Google DeepMind scholarship for funding my PhD.

I thank my dear friends for always lending an ear and giving me words of encouragement and support. They are always a great comfort during challenging times.

Finally, I am eternally grateful to my family for their constant love and support. My love and thanks to my parents, my brother, and our dog, Patrick.

Abstract

The advent of deep learning has brought about great progress on many fundamental computer vision tasks such as classification, detection, and segmentation, which describe the categories and locations of objects in images and video. There has also been much work done on supervised learning—teaching machines to solve these tasks using human-annotated labels. However, it is insufficient for machines to know only the names and locations of certain objects; many tasks require a deeper understanding of the complex physical world—how objects interact with their surroundings, for example (often by creating shadows, reflections, surface deformations, and other visual effects). Furthermore, training models to solve these tasks while relying heavily on human supervision is costly and impractical to scale. Thus, this thesis explores two directions: first, we aim to go beyond segmentation and address a wholly new task: grouping objects with their correlated visual effects (e.g. shadows, reflections, or attached objects); second, we address the fundamental task of video object segmentation in a self-supervised manner, without relying on any human annotation.

To automatically group objects with their correlated visual effects, we adopt a layered approach: we aim to decompose a video into object-specific layers which contain all elements moving with the object. One application of these layers is that they can be recombined in new ways to produce a highly realistic, altered version of the original video (e.g. removing or duplicating objects, or changing the timing of their motions). Here the key is to leverage natural properties of convolutional neural networks to obtain a layered decomposition of the input video. We design a neural network that outputs layers for a video by overfitting to the video. We first introduce a human-specific method, then show how it can be adapted to arbitrary object classes, such as animals or cars.

Our second task is video object segmentation: producing pixel-wise labels (segments) for objects in videos. Whereas our previous work is optimized on a single video, here we take a data-driven approach and train on a large corpus of videos in a self-supervised manner. We consider two different task settings: (1) semi-supervised object segmentation, where an initial object mask is provided for a single frame and the method must propagate this mask to the remaining frames, and (2) moving object discovery, where no mask is given and the method must segment the salient moving object. We explore two different input streams: RGB and optical flow, and discuss their connection to the human visual system.

Contents

List of Abbreviations	ix
1 Introduction	1
1.1 An Application: Video Editing	2
1.2 Contributions and Outline	4
1.3 Publications	7
2 Background Literature	9
2.1 Layered Video Representations	9
2.1.1 Formulation	9
2.1.2 Layer Decomposition Methods	11
2.2 Video Completion	14
2.3 Neural Rendering	14
2.4 Self-Supervised Learning from Videos	15
2.4.1 Recognition and Retrieval	16
2.4.2 Audiovisual Tasks	18
2.4.3 Object Structure	19
2.4.4 Object Segmentation	19
3 Layered Neural Rendering for Retiming People in Video	23
4 Omnimatte: Associating Objects and Their Effects in Video	39
5 MAST: A Memory-Augmented Self-Supervised Tracker	61
6 Self-Supervised Video Object Segmentation by Motion Grouping	77
7 Class-Agnostic Counting	97
8 Conclusion	115
8.1 Achievements and Impact	115
8.2 Future Work	118
8.2.1 Faster Omnimatte	118
8.2.2 End-to-End Training	119

List of Abbreviations

CNN	Convolutional Neural Network
DIP	Deep Image Prior
VOS	Video Object Segmentation

Stop playing and help me find my shadow!

— Peter Pan

1

Introduction

In this thesis we aim to design computer vision algorithms that can understand how objects in videos affect their environment. As people, animals, and objects move through the world, they cast shadows, reflections, create ripples in water, deform surfaces, etc. When observing this phenomena in the world, we are able to immediately understand which of these *effects* is caused by which moving entity. We even have the ability to imagine what the scene would look like if that entity never existed—those same effects would also never appear.

Existing computer vision techniques cannot automatically group objects with their correlated visual effects. While the field has seen much progress in segmenting and tracking objects in videos, effects such as shadows, reflections, and surface deformations are typically ignored or handled by a dedicated system (e.g. a shadow detector) rather than a generalized solution. Understanding which scene elements in a video are related to one another is a fundamental reasoning task that we explore in this thesis.

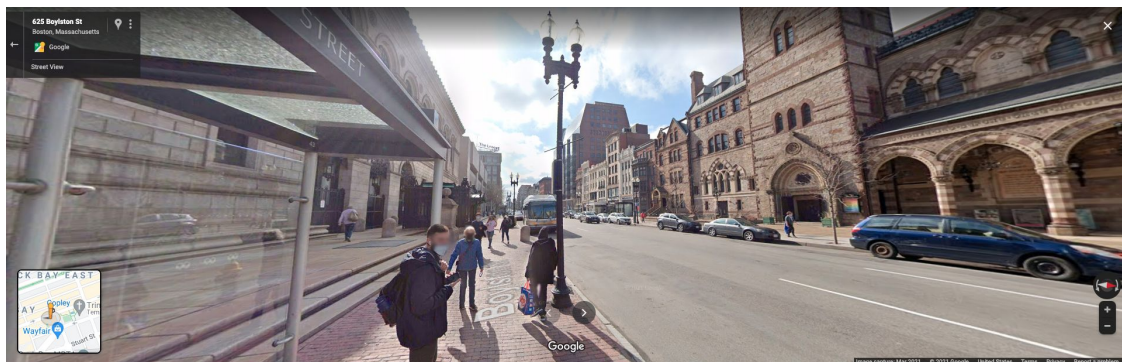
Our goal is to automatically infer from a video the effects that are correlated with a moving object. Furthermore, we wish to do this *unsupervised*; that is, without requiring any human annotation, by merely observing the raw input video. The signal we will leverage to learn this association is *motion*—objects and their effects move in correlated ways.

In the first part of this thesis, we present methods to automatically group objects with their correlated effects by leveraging natural properties of deep convolutional neural networks. For this work, we first use existing off-the-shelf tools to detect and segment the moving objects whose effects we wish to label. In the second part of the thesis we investigate ways to improve the aforementioned tools for video object segmentation. Our focus is once again to learn *without* any human annotations, by taking advantage of natural signals present in raw video: e.g. multiple viewpoints, smooth and gradual changes, all of which enable the learning of *correspondence* in a fully self-supervised manner.

1.1 An Application: Video Editing

While the goal of this thesis is to *understand* video, one downstream application of our work is that it can also be used to *edit* and *re-render* visual content. Editing images and videos has a wide range of uses, among which include the removal of individuals to maintain their privacy, stroboscopic images for sports visualization and analysis, and creative expression (see Figure 1.1). Creating such edits where individual objects in a scene are manipulated independently of their surroundings is extremely challenging and costly, as it requires *object mattes*, which contain pixel-level information about the image regions that belong to the object. To acquire these mattes, either the surroundings must be strictly controlled and staged (e.g. with the use of greenscreen), or when this setup is not possible, the object matte is often drawn by hand through a process called rotoscoping. The costliness of this process makes high-quality editing largely confined to the professional film industry and highly skilled individuals rather than casual users.

Fortunately, computer vision has already made strides in automating aspects of these tasks, saving users much time and effort. Automatic image matting (and video matting, to a lesser extent) is an active area of research. However, these methods are all concerned with identifying solely the object, and do not take into account any correlated visual effects such as shadows or reflections. Ignoring these elements, however, results in unrealistic scenes where people are



(a)



(b)



(c)

Figure 1.1: Sample applications of image and video editing. (a) A screenshot taken of Google Streetview, where faces are blurred to preserve the privacy of individuals. (b) A composite image created by fusing different video frames to visualize a goal scored during a soccer match. (c) A sequence from the film *Chungking Express* [Wong 1994] which uses different speeds of the people in the scene for artistic effect. The person in the foreground moves in slow motion while the people behind him rush by in a quick blur.

removed or edited, but their shadows and reflections remain unchanged. While there are dedicated shadow detectors, even more complex scenes may contain water splashes or surface deformations, and no existing method can capture all these different types of visual effects.

1.2 Contributions and Outline

In this section, we provide a summary of the main contributions of this thesis, and outline the chapters.

Grouping people with their correlated effects. In Chapter 3, we propose the new task of grouping people in video with their *correlated visual effects* (e.g. shadows, reflections, attached objects), and address this challenge with a novel video layer decomposition method based on neural rendering [Thies, Zollhöfer, and Nießner 2019]. The human representation that we use is UV maps. To obtain these maps, we first use off-the-shelf tools to track the people in the input video, and obtain their keypoints. Next, we use our novel *keypoint-to-UV network*, which we pretrained offline, to convert the keypoints to UV maps. We input the UV map for each person to our convolutional neural network, which predicts the RGBA layer for the corresponding person. The RGBA layers are then back-to-front composited to reconstruct the original video frame. The network is optimized on a single video and is trained to minimize reconstruction loss. Once trained, these layers can be used to produce a variety of video editing effects, such as person removal, duplication, or *retiming*—changing the speed of individual people. We demonstrate results on videos of complex scenes, including elements such as water, trampolines, and loose clothing.

Grouping objects with their correlated effects. In Chapter 4, we extend the human-specific method presented in Chapter 3 to operate on objects of *any* class. We show that the UV geometry inputs can be replaced by binary object masks and optical flow to provide frame-to-frame correspondence. We additionally add flow-based losses for improved layer separation and temporal consistency, which we demonstrate in our ablations. Furthermore, while our original work in Chapter 3 represents the camera motion with a homography, we additionally optimize for a coarse background warp to allow for greater expressivity, and show that this produces cleaner object layers. We demonstrate our method qualitatively on challenging videos containing diverse subjects (e.g. animals, cars, people) and

diverse effects (e.g. shadows, smoke, reflections). We additionally evaluate our method quantitatively on data labeled for change detection. We also qualitatively compare our approach with methods specific to object removal [Gao et al. 2020] and shadow detection [Tianyu Wang et al. 2020] and show comparable if not favorable results, despite ours being a general approach.

Self-supervised video object segmentation. In Chapter 5, we present a method for tracking object segments in video. Our method is trained in a fully self-supervised manner, i.e. without relying on any human annotation. We investigate various training settings that have been used in prior self-supervised object segmentation works and conduct thorough ablation studies to determine the optimal choices. We additionally incorporate a memory module to overcome failures due to severe occlusion. At the time of publication, our method surpassed the state-of-the-art self-supervised methods on existing benchmarks, while approaching the performance of supervised methods. We additionally show that self-supervised methods have a lower *generalization gap* between seen and unseen object classes than supervised methods, and argue that this metric should be considered as it reflects real-world use cases.

Self-supervised moving object discovery. In Chapter 6, we present a self-supervised method for discovering and segmenting moving objects in videos. Our method exploits motion cues by operating only on precomputed optical flow [Teed and J. Deng 2020], in contrast to previous methods that perform segmentation directly on RGB inputs. We adopt a simple variant of the Transformer [Vaswani et al. 2017] to decompose the input flow into a foreground object layer and background layer, and train the network in a self-supervised manner by re-compositing the layers to reconstruct the original flow. At the time of publication, our approach surpassed previous state-of-the-art self-supervised methods on public benchmarks, while also being an order of magnitude faster. We additionally evaluate our method on the challenging task of segmenting camouflaged objects and significantly

outperform other self-supervised approaches while being competitive with the top supervised approach.

Class-agnostic counting. In Chapter 7, we take a brief detour from video understanding to consider the problem of counting objects of arbitrary classes. Existing counting methods are designed for a specific class of object; however, we observe that humans are capable of counting objects despite having never seen them prior (see Figure 1.2). Thus, we propose to reformulate counting as a *matching* problem: given an ‘exemplar’ image patch containing an instance of the object to count (e.g. a car, a cell), and an image to count, our method outputs a heatmap that localizes all occurrences of the object. The heatmap supports both counting by detection (via taking local maximums) and counting by density estimation (by summing the heatmap), the latter of which is the preferred method in the case of heavily overlapping object instances. Returning to this thesis’ theme of *videos*, we first pretrain our generic counting-by-matching network on Imagenet Video data [Russakovsky et al. 2015] labeled for object tracking; this allows the network to learn to match different views of the same object. This pretrained network can then be easily adapted to new domains and user requirements by finetuning just 6% of network parameters on small amounts of data. We demonstrate our method on existing benchmarks for synthetic cells, real cells, cars, and human crowds, and achieve performance at or surpassing state-of-the-art methods.



Figure 1.2: Counting unseen object classes. Humans have the ability to count objects even if they have never seen them before, due to our knowledge of intra-class variations in color and shape, for example. Even if a person has never seen this type of fruit before, they would have no trouble counting it, despite individual differences in appearance.

1.3 Publications

This thesis consists of the following published works:

- Chapter 3: [Erika Lu](#), Forrester Cole, Tali Dekel, Weidi Xie, Andrew Zisserman, David Salesin, William T. Freeman, and Michael Rubinstein. Layered Neural Rendering for Retiming People in Video. In *SIGGRAPH Asia, 2020*.
- Chapter 4: [Erika Lu](#), Forrester Cole, Tali Dekel, Andrew Zisserman, William T. Freeman, and Michael Rubinstein. Omnimatte: Associating Objects and Their Effects in Video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2021*.
- Chapter 5: Zihang Lai, [Erika Lu](#), and Weidi Xie. MAST: A Memory-Augmented Self-Supervised Tracker. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020*.

- Chapter 6: Charig Yang, Hala Lamdouar, Erika Lu, Andrew Zisserman, and Weidi Xie. Self-supervised Video Object Segmentation by Motion Grouping. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV), 2021*.
- Chapter 7: Erika Lu, Weidi Xie, and Andrew Zisserman. Class-Agnostic Counting. In *Asian Conference on Computer Vision (ACCV) 2018*.

Note on formatting. This thesis is presented as an *integrated thesis*¹, where publications are presented in the original format in which they were published.

¹<https://academic.admin.ox.ac.uk/annex-b-integrated-theses-guidance-for-divisional-boards>

2

Background Literature

This thesis has two primary focuses: (1) decomposing videos into layers of objects with their correlated effects, and (2) segmenting objects in videos, both of which are done without any manual annotation. Here we discuss the relevant background literature. Section 2.1 first introduces layered representations before providing an overview of layer decomposition methods with a focus on deep-learning based approaches. Section 2.3 discusses neural rendering. In Section 2.4, we cover self-supervised learning from videos, before delving deeper into self-supervised methods specific to the video object segmentation task (Section 2.4.4).

2.1 Layered Video Representations

2.1.1 Formulation

Early works to propose a layered representation for video include Adelson [1991](#); Darrell and Pentland [1991](#); J. A. Wang and Adelson [1993a](#); J. A. Wang and Adelson [1993b](#); J. A. Wang and Adelson [1994](#), with the idea being that a video may contain many complex motions, but it can be explained by individual layers which each contain simpler motion (see Figure 2.1). Formally, a frame I can be decomposed into a set of RGBA (color and opacity) layers:

$$\mathcal{L} = \{L_i\}_{i=0}^{N-1} = \{E_i, \alpha_i\}_{i=0}^{N-1}, \quad (2.1)$$

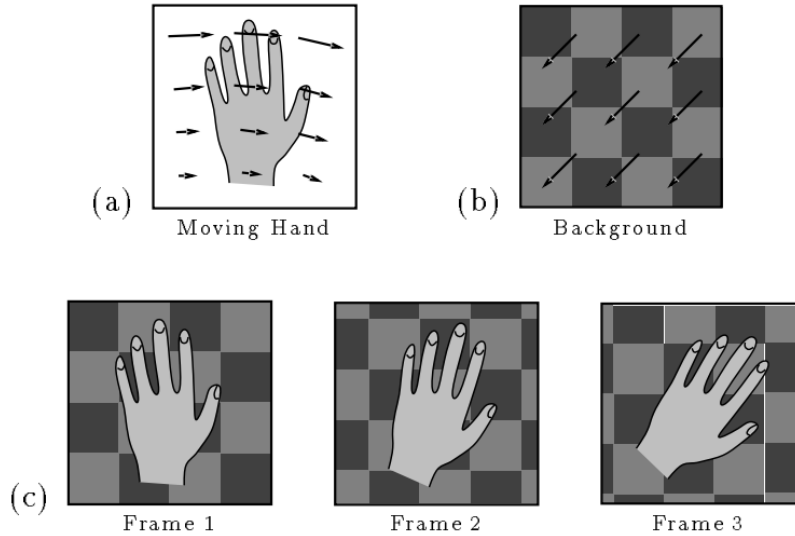


Figure 2.1: Representing a video as layers. A figure from J. A. Wang and Adelson 1994, the first work to propose a layered representation for video. (a) A moving hand, (b) the background, translating downwards and to the left, and (c) the observed image sequence. Reproduced from J. A. Wang and Adelson 1994.

where N is the number of layers; E_i is the color map for layer i ; and α_i is the alpha map for layer i , which controls the opacity of the layer and takes a value in the range $[0, 1]$ to deal with transparency, shadows, motion blur, etc. The layers are then back-to-front composited to reconstruct the original video frame.

To produce a sequence of video frames, the layers are warped according to the *velocity map* V . For greater representational flexibility, a *delta map* D may additionally update the color map over time to account for changes that cannot be captured by the layered model. The full layer compositing model that includes velocity maps and delta maps can be seen in Figure 2.2.

The advantage of the layered model is that new image sequences can be easily rendered by combining the layers in new ways. The challenge in decomposing a sequence of images into layers lies in jointly estimating the assignment of pixels to layers and the motions of the layers. To quote J. A. Wang and Adelson 1994, “Synthesis is easy but analysis is hard.” J. A. Wang and Adelson 1994 use a two-stage approach to estimating layers: 1) motion estimation, and 2) segmentation by affine model fitting.

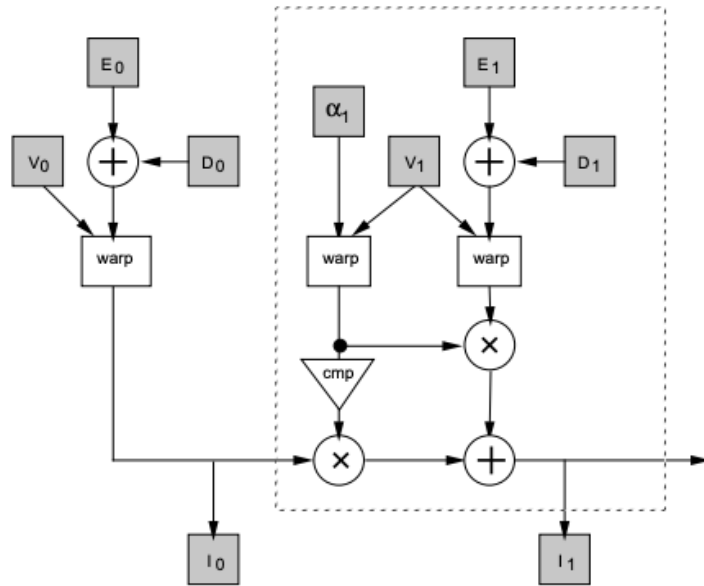


Figure 2.2: Diagram of layer compositing proposed by J. A. Wang and Adelson 1994. Each layer E_i is updated by a delta map D_i and warped by a velocity map V_i . Transparency map α_i determines how to blend current layer E_i with output I_{i-1} from the previous compositing step. Reproduced from J. A. Wang and Adelson 1994

2.1.2 Layer Decomposition Methods

Layered models have since been applied to a range of tasks, including view synthesis and reflection removal. Many works have adapted the layered model definition to represent depth layers rather than motion layers, enabling novel view synthesis [Baker et al. 1998; Shade et al. 1998; Torr et al. 2001; Zitnick et al. 2004]. Brostow and Essa 1999 use temporal occlusions to estimate the relative depths of layers. Jojic and Frey 2001 extend the layered model to introduce the *flexible sprite model*, which probabilistically model moving, occluding objects. Kumar et al. 2008 extend this model to include the effects of lighting and motion blur. Fradet et al. 2008 present an algorithm for motion segmentation and mosaic generation that can accommodate human-in-the-loop correction. Layered models have also been used for the goal of estimating optical flow, as the model is able to handle the sharp discontinuities at surface boundaries [D. Sun, Sudderth, et al. 2012; D. Sun, Wulff, et al. 2013; Wulff and Black 2014; Wulff and Black 2015; Sevilla-Lara et al. 2016]. Wulff and Black 2014 additionally use a layered representation to model motion

blur, enabling them to perform deblurring. Xue et al. 2015; Nandoriya et al. 2017 use a layered model to remove reflections and obstructions by taking advantage of the different motions of the layers due to visual parallax.

Layer Decomposition With Deep Neural Networks. More recently, several works have leveraged deep neural networks to perform layer decomposition, taking advantage of deep networks’ ability to learn from large amounts of data [Alayrac, Carreira, and Zisserman 2019; Alayrac, Carreira, Arandjelovic, et al. 2019], as well as properties arising from the structure of convolutional neural networks [Gandelsman et al. 2019]. T. Zhou et al. 2018 train a neural network on a pair of stereo images to predict layers that correspond to physically accurate depth, which they refer to as *multiplane images* (MPI). They use video clips of static scenes to train the network to produce MPIs from narrow-baseline stereo image pairs. The predicted layers contain disoccluded content, which enables novel-view synthesis. Srinivasan, Tucker, et al. 2019 build upon this work and achieve significantly larger view extrapolations. Flynn et al. 2019 extend the original MPI work by applying learned gradient descent [Adler and Oktem 2017; Adler and Oktem 2018], allowing it to require much fewer update steps at inference time than in the previous work. Tucker and Snavely 2020 build on this line of work to predict MPIs from a *single* image.

Alayrac, Carreira, and Zisserman 2019 use a learning-based approach to separate videos into a fixed number of layers, for applications such as reflection and haze removal. The authors generate synthetic training data by sampling two random videos and blending them with a random mixing parameter $\alpha \in [0, 1]$. The videos are taken from the Kinetics-600 dataset [Kay et al. 2017], which is comprised of humans performing a variety of actions. Despite training on synthetic data, they demonstrate results on real videos that include shadows, reflections, and smoke. Alayrac, Carreira, Arandjelovic, et al. 2019 extend this work by making the layer decomposition process controllable: the layers are conditioned on external cues such as audio.

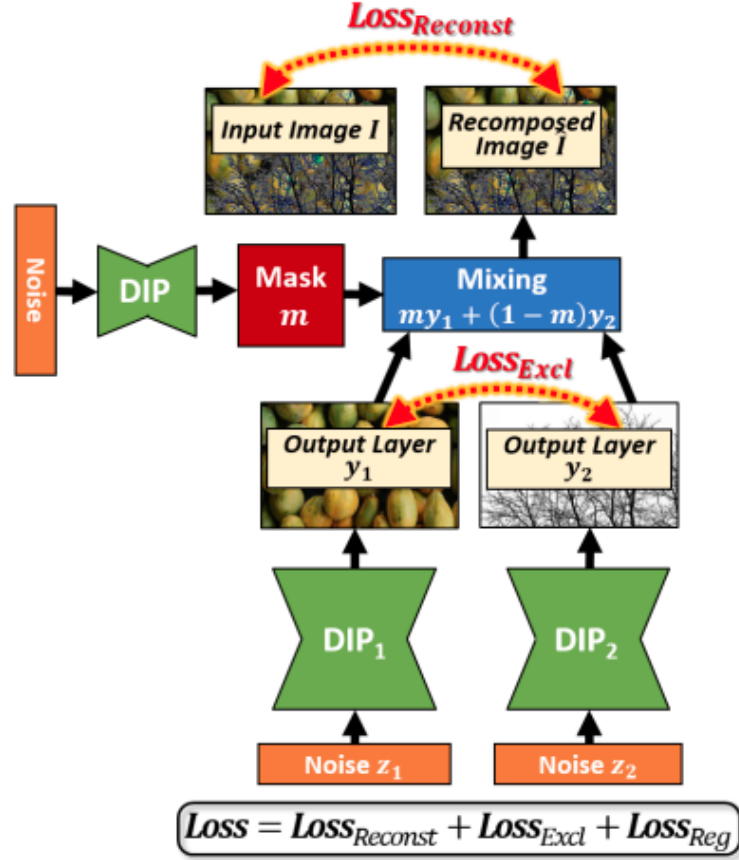


Figure 2.3: Double-DIP method overview. Two DIP networks each predict an output layer y from fixed input noise z . A third DIP network predicts a mask m which is used to mix the output layers y_1 and y_2 to reconstruct the original image. The DIP networks are optimized to minimize reconstruction loss, an exclusion loss on y_1 and y_2 , and a task-specific regularization loss. Reproduced from Gandelsman et al. 2019.

Double-DIP [Gandelsman et al. 2019] is a method that leverages the “deep image prior” (DIP) property [Ulyanov et al. 2018] of convolutional neural networks to perform layer separation. The DIP property refers to a finding in Ulyanov et al. 2018 that the *structure* of CNNs gives them the tendency to output smooth, ‘natural’ images rather than random noise images. They demonstrate this behavior by training a single CNN to overfit to a real image which is corrupted by adding random noise, JPEG artifacts, or arbitrary masking. During training, the network first outputs a ‘clean’ version of the image *without* the corruption, before eventually overfitting to the artifacts, thus demonstrating CNNs’ bias towards producing realistic images. The Double-DIP method utilizes this property to produce two

clean layers from a single image. Two neural networks with identical architectures are each passed an input of fixed random noise, and optimized to output layers that when composited, reconstruct the original image. Since the networks are initialized with different weights and given different noise inputs, they each capture a different layer of the image. They demonstrate results on foreground/background segmentation, dehazing, transparency separation, and even watermark removal. While the method is primarily for images, the authors demonstrate its effectiveness on videos as well, by initializing different random noise for each frame.

2.2 Video Completion

Video completion, or video inpainting, is the task of filling in missing pixels in a video. Applications include restoration (e.g. fixing scratches), watermark removal, and object removal. While object removal is also one application of our work in Chapters 3 and 4, we can think of video completion as *complementary* to our layer decomposition methods. Video completion takes as input a mask indicating the pixels for which to synthesize content. In Chapters 3 and 4, we present methods for producing a matte for a given video, which marks the pixels belonging to an object and the ‘effects’ created by that object (e.g. shadows, reflections). This matte can be converted to a mask via thresholding and used as input to a video completion method to completely remove the object’s presence (effects included).

Traditional methods for video completion typically use patch-based synthesis techniques (e.g. [Wexler et al. 2007, Newson et al. 2014]). A limitation of these techniques is that they reuse existing patches in the video, and typically have high processing time. Several works use optical flow to guide patch synthesis (e.g. [Huang et al. 2016]), or to directly inpaint by color propagation [Bokov and Vatolin 2018; Okabe et al. 2020]. Deep learning-based methods have seen some success with 3D convolutions [Chang et al. 2019; C. Wang et al. 2019]; however, they are limited in resolution due to the high memory costs of 3D convolutions. Kim, Woo, et al. 2019; S. Lee et al. 2019; Oh et al. 2019 aim to address this issue by using nearby frames

as references. Recently, several methods have seen success in combining flow-based methods with deep learning [R. Xu et al. 2019; Gao et al. 2020].

2.3 Neural Rendering

Neural rendering is a new and rapidly growing field that leverages neural networks for rendering photorealistic imagery [Tewari et al. 2020]. Traditional computer graphics uses hand-crafted scene representations and sophisticated graphics engines to produce high-quality images. However, generating scene elements such as lighting, materials, and shape is difficult to do automatically. Neural rendering provides an attractive alternative where the complexities of scene representations and rendering are abstracted away from the user and handled implicitly by the neural network, at the expense of some control. Neural rendering methods combine deep generative models with traditional computer graphics techniques. The goal of many such systems is *generalization*: once trained on a dataset of examples, the system should be able to render new photorealistic content (a new configuration of previously observed scene properties). Neural rendering has been used for numerous applications such as novel view synthesis [Hedman et al. 2018; Nguyen-Phuoc et al. 2018; Meshry et al. 2019; Sitzmann, Thies, et al. 2019; Sitzmann, Zollhöfer, et al. 2019; Thies, Zollhöfer, Theobalt, et al. 2020; Mildenhall et al. 2020; Martin-Brualla et al. 2021], facial and body reenactment [Lombardi et al. 2018; Thies, Zollhöfer, and Nießner 2019; Chan et al. 2019; Liu et al. 2019], semantic photo manipulation [J. Zhu et al. 2016; Q. Chen and Koltun 2017; Isola, J.-Y. Zhu, et al. 2017; Ting-Chun Wang et al. 2018; Bau et al. 2019; Park et al. 2019], and relighting [Z. Xu et al. 2018; H. Zhou et al. 2019; T. Sun et al. 2019; Zhang et al. 2021; Srinivasan, B. Deng, et al. 2021].

While generalization is a goal of neural rendering, our work in Chapter 3 differs in that generalization to new scenes is not the aim; rather, our work employs a neural rendering architecture as a *building block* for decomposing a video into layers. Our method builds upon a similar framework used in Thies, Zollhöfer, and Nießner 2019; Aliev et al. 2020. The input scene geometry is represented by UV

texture coordinates [Thies, Zollhöfer, and Nießner 2019] or a point cloud [Aliev et al. 2020]. The geometry is then mapped to a learnable ‘neural texture’ map before being passed to a rendering network, which outputs a photorealistic image. During inference time, the input geometry can be manipulated by the user to produce new scenes. Aliev et al. 2020; Thies, Zollhöfer, and Nießner 2019 demonstrate results on view synthesis; Thies, Zollhöfer, and Nießner 2019 additionally apply their method to scene editing and facial reenactment.

2.4 Self-Supervised Learning from Videos

Due to the high cost of manual annotation and the increase in available unlabeled video data, there has been much interest in learning useful feature representations *without* requiring any human supervision. Self-supervised learning typically involves constructing a *proxy task* which the network can only solve by learning a meaningful feature representation. These features can then be used for downstream tasks such as video retrieval, action recognition, segmentation, or tracking.

2.4.1 Recognition and Retrieval

A variety of proxy tasks have been proposed for learning representations at the frame-level that are useful for video retrieval or human action recognition. Standard benchmarks used for evaluating human action recognition are HMDB51 [Kuehne et al. 2011], UCF101 [Soomro et al. 2012], and Kinetics [Kay et al. 2017]. Two common evaluation settings are: (1) linear probing, where the learned features are frozen, and a single linear layer is trained with cross-entropy loss, and (2) finetuning, where the feature encoder and the linear layer are both finetuned.

Misra et al. 2016; Fernando et al. 2017; H.-Y. Lee et al. 2017; Wei et al. 2018; D. Xu et al. 2019; Kim, Cho, et al. 2018 use the natural spatio-temporal ordering of patches or frames as a signal for learning a semantically meaningful representation. Misra et al. 2016 sample video clips that are either in order or shuffled (see Figure 2.4a), and then train a network to classify whether the input sequence of frames is temporally correct (Figure 2.4b). Misra et al. 2016 show that

their trained features are more sensitive to human pose than ImageNet [J. Deng et al. 2009] pretrained or randomly initialized network features (Figure 2.4c). Thus, the self-supervised temporal order verification of Misra et al. 2016 is a more suitable pretraining task for human action recognition than is fully supervised ImageNet classification, which is the method that had been typically used in the literature.

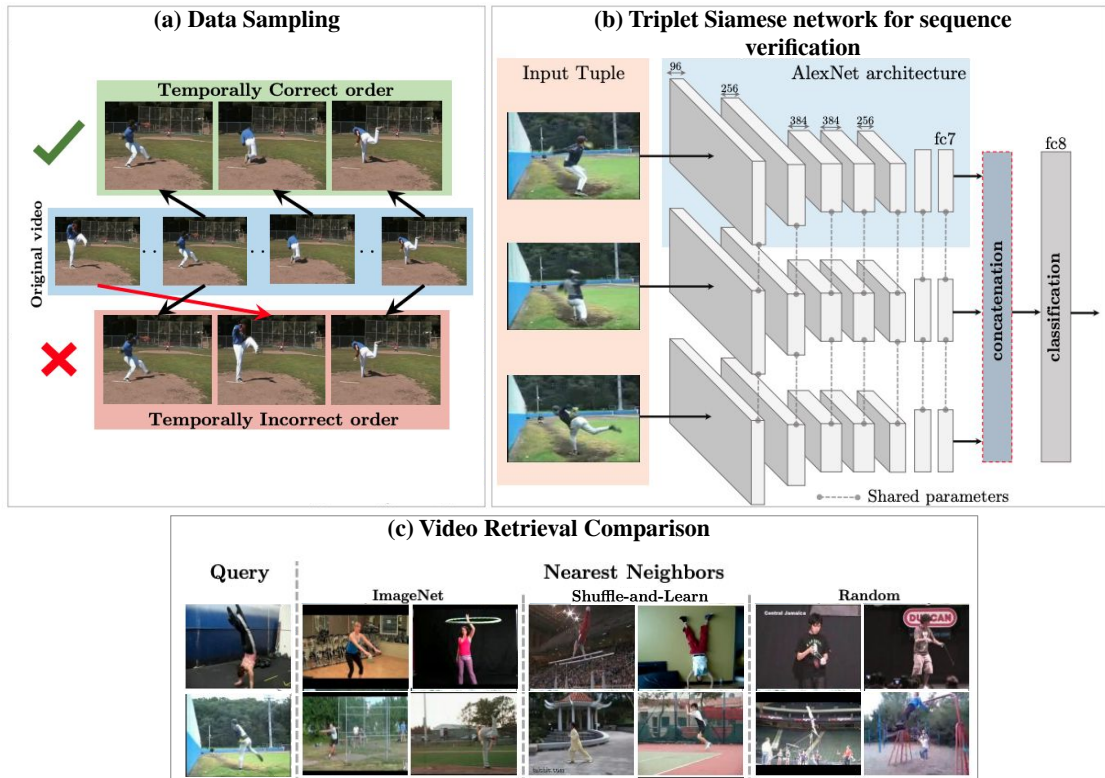


Figure 2.4: Method and Results for ‘Shuffle-and-Learn’ [Misra et al. 2016]. (a) Data is sampled from videos by arbitrarily shuffling or preserving the original temporal order. (b) The sequence of frames is passed to a network, which is trained to predict whether the data is in the correct temporal order. (c) For each query video, nearest neighbor videos are retrieved using ImageNet pretrained features, the Shuffle-and-Learn features, and randomly initialized network features. In comparison to ImageNet and random, Shuffle-and-Learn features are more sensitive to human pose, returning nearest neighbors that contain the same pose, while ImageNet nearest neighbors are of a similar image class (top: exercise, bottom: baseball). Reproduced with modification from Misra et al. 2016.

In addition to spatiotemporal ordering, speed prediction has also been used as a proxy task [Benaim et al. 2020; Epstein et al. 2020; J. Wang et al. 2020]. Benaim et al. 2020 train a network to predict whether clips have been randomly sped up or slowed down. They demonstrate that the resulting features can be used to select

temporal sequences that appear similar regardless of changes in speed, enabling videos to be ‘intelligently’ sped up in a manner that is less perceptually disruptive to users. Epstein et al. 2020 address a new task of predicting unintended actions in videos and find speed prediction to be the most effective proxy task.

Future frame prediction and generation has also been extensively explored for learning useful feature representations [Vondrick, Pirsiaavash, et al. 2016a; Vondrick, Pirsiaavash, et al. 2016b; Liang et al. 2017; Vondrick and Torralba 2017; Denton and Birodkar 2017; Diba et al. 2019; Han et al. 2019; Han et al. 2020a]. Han et al. 2020a; Han et al. 2020b combine contrastive learning with future frame prediction to learn representations for action recognition and video retrieval. Such contrastive methods learn to discriminate a transformed video sample against other samples in the dataset. The transformation can be synthetic (such as data augmentation) or natural (such as two different clips within the same video).

Other signals that have been leveraged for self-supervised representation learning include motion [Agrawal et al. 2015; Jayaraman and Grauman 2015; Jayaraman and Grauman 2016], geometry cues [Gan et al. 2018], and rotation [Jing et al. 2018]. Isola, Zoran, et al. 2015 use co-occurrence of visual entities at various levels: image patches, for image segmentation; frames, for movie scene segmentation; and geospatial proximity in image collections, for clustering. X. Wang and Gupta 2015 extract and track patches in videos and use positive and negative pairs to train a Siamese Triplet Network.

2.4.2 Audiovisual Tasks

There is a growing body of work on taking advantage of multiple modalities such as the *audiovisual* to learn self-supervised representations. Alwassel et al. 2019; Arandjelović and Zisserman 2017; Arandjelović and Zisserman 2018; Korbar et al. 2018; Patrick et al. 2020; Piergiovanni et al. 2020; Afouras, Asano, et al. 2021 use a contrastive loss to learn the correspondence between video frames and the accompanying audio. Narrated video has also been leveraged by Miech et al. 2020

to train joint video and text embeddings, enabling text-to-video retrieval in addition to action recognition and localization.

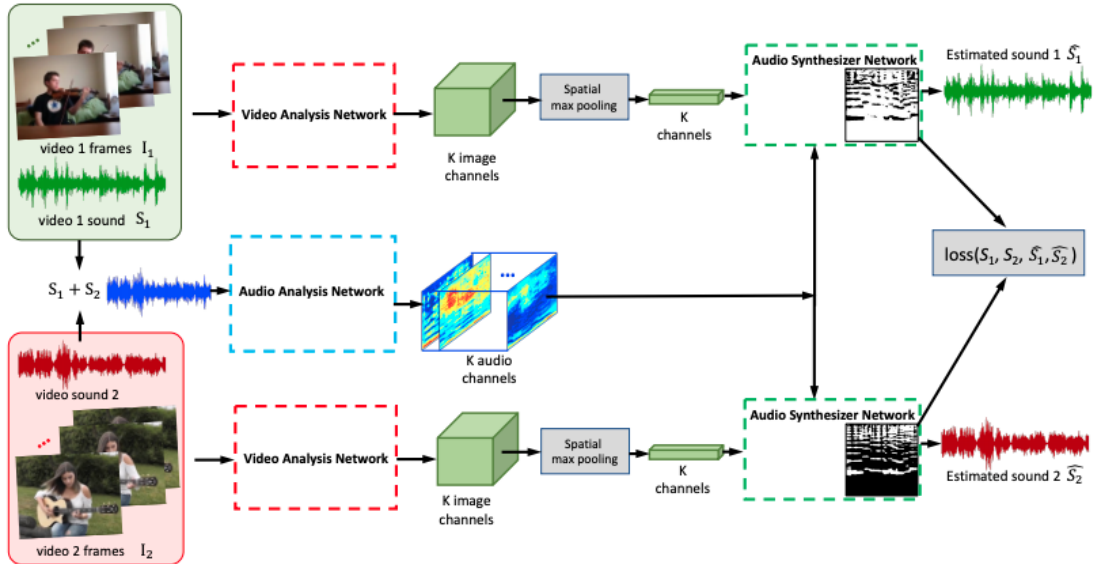


Figure 2.5: Method overview for self-supervised audiovisual source separation [Zhao, Gan, Rouditchenko, et al. 2018]. The audio tracks from two different videos are blended together. The blended audio and each video’s frames are passed to the audio analysis network and video analysis network, respectively. The outputs are passed to the audio synthesizer network, which takes in the audio features and a single video’s features, and predicts only the audio that corresponds to that input video. Reproduced from Zhao, Gan, Rouditchenko, et al. 2018.

Multi-modal representation learning has applications in audiovisual source separation and localization [Zhao, Gan, Rouditchenko, et al. 2018; Zhao, Gan, Ma, et al. 2019; Harwath et al. 2018; Afouras, Chung, et al. 2018b; Arandjelović and Zisserman 2018; Afouras, Chung, et al. 2020], lip reading [Afouras, Chung, et al. 2018a; Chung and Zisserman 2016; Chung, Senior, et al. 2017], and audiovisual navigation [C. Chen et al. 2020]. Zhao, Gan, Rouditchenko, et al. 2018 train a self-supervised model to separate sounds conditioned on image pixels. They mix the audio tracks from two different videos and train a model to separate the blended audio by conditioning on the visual stream from each video (Figure 2.5). The authors demonstrate the model’s ability to separate the sounds produced by different musical instruments, suggesting that it has learned semantics despite not having access to class labels.

2.4.3 Object Structure

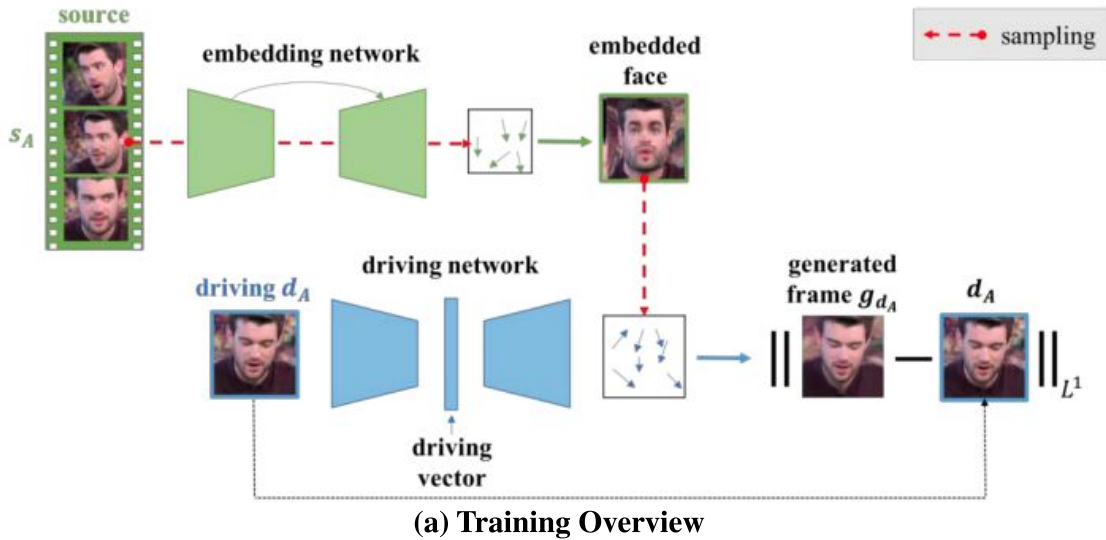


Figure 2.6: Overview of self-supervised face synthesis [Wiles et al. 2018b]. (a) During training, an embedding network predicts an embedded face (encoding person identity in a neutral pose) from the source frames. The driving network takes as input an image containing the desired pose, and predicts a warp to apply to the embedded face. The networks are optimized to match the synthesized face to the original driving image. (b) In comparison to CycleGAN [J.-Y. Zhu et al. 2017], X2Face [Wiles et al. 2018b] better disentangles identity from pose. Reproduced from Wiles et al. 2018b.

Several recent deep-learning based self-supervised works are targeted towards learning object structure from unlabeled videos. Wiles et al. 2018b; Wiles et al. 2018a train a network to warp a source video frame of a face to a target frame from the same face track, and demonstrate that the learned embeddings disentangle identity from pose and expression. (Figure 2.6). Jakab et al. 2018; Jakab et al. 2020 impose an information bottleneck to learn human pose and object landmarks.

Wu et al. 2021 trains an object category-specific model on videos to predict a textured, articulated 3D mesh from a single image.

2.4.4 Object Segmentation

The self-supervised work in chapters 5 and 6 of this thesis address the video object segmentation task. The goal of video object segmentation (VOS) is to output a pixel-wise binary mask for each frame that indicates the presence of a specified object. There are two common test settings: (1) *semi-supervised VOS*, where an initial object mask is given for a frame, and the object must be tracked in the remaining frames, and (2) *zero-shot VOS*, where no object mask is given, and the object (typically the most salient) must be discovered and segmented. Despite this terminology, the *training* setting for both tasks can be either supervised or unsupervised (e.g. the semi-supervised VOS task can be trained in an unsupervised manner). Zero-shot VOS with a focus on moving objects can also be referred to as moving object discovery, or motion segmentation. Here we discuss self-supervised training methods for semi-supervised VOS, with a focus on deep-learning based methods.

Tracking by colorizing. Our work on self-supervised training for object segmentation follows a similar framework laid out in Vondrick, Shrivastava, et al. 2018. Vondrick, Shrivastava, et al. 2018 leverage the natural temporal consistency of color in videos by using *colorization* as a proxy task for learning to track. Given two grayscale images I_1 and I_2 from a video, the model is trained to colorize each pixel in I_2 by *copying* colors from pixels in I_1 (see Figure 2.7). This training regime encourages the model to naturally learn correspondence between video frames, by optimizing for the CNN parameters that result in similar pixels being embedded in the same space. During training, the colors are copied from the reference frame to the target frame, whereas during inference time, the class labels are propagated instead in order to solve the segmentation task.

A weakness of Vondrick, Shrivastava, et al. 2018 is that the inputs to the network must be grayscale; thus, the trained model can never make use of color information,

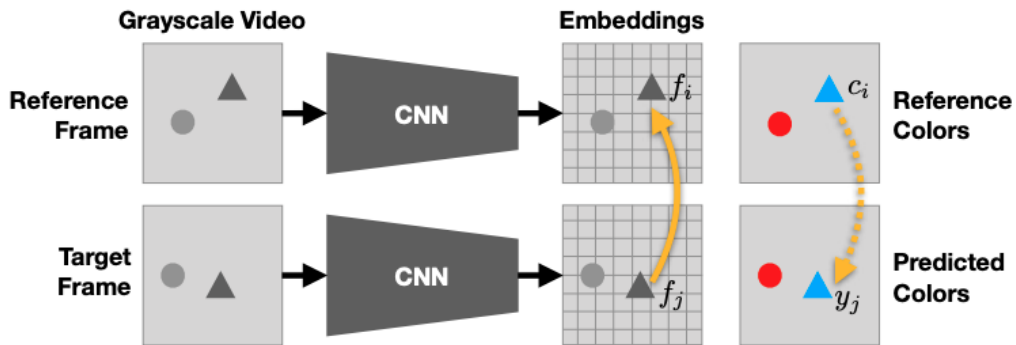


Figure 2.7: Method overview for learning to track by coloring. Two grayscale video frames are each embedded by a CNN. Each pixel in the target frame is then assigned a color value by ‘copying’ the colors from the most similar pixels in the reference frame (based on softmax similarity). For visualization purposes, the copying/tracking operation is drawn with a single yellow arrow, but in practice this tracking is soft. Reproduced from Vondrick, Shrivastava, et al. 2018.

lacking a strong cue for the tracking task. Follow-up work by Lai and Xie 2019 addresses this issue by passing as inputs the original RGB video frames with color jittering and augmentations applied. They additionally incorporate scheduled sampling and cycle consistency for improved performance on the DAVIS-2017 [Pont-Tuset et al. 2017] benchmark. X. Wang, A. Jabri, et al. 2019 rely primarily on a cycle consistency loss to train a self-supervised network to track patches forwards and backwards. Lu et al. 2020 also use the cycle consistency principle in a framework that accommodates both unsupervised and weakly supervised training, as well as zero-shot and semi-supervised VOS.

3

Layered Neural Rendering for Retiming People in Video

This work was presented at the 13th ACM SIGGRAPH Conference and Exhibition on Computer Graphics & Interactive Techniques in Asia, 2020.

Layered Neural Rendering for Retiming People in Video

ERIKA LU, University of Oxford, Google Research

FORRESTER COLE, Google Research

TALI DEKEL, Google Research

WEIDI XIE, University of Oxford

ANDREW ZISSERMAN, University of Oxford

DAVID SALESIN, Google Research

WILLIAM T. FREEMAN, Google Research

MICHAEL RUBINSTEIN, Google Research

Input Video



Our Retiming Result

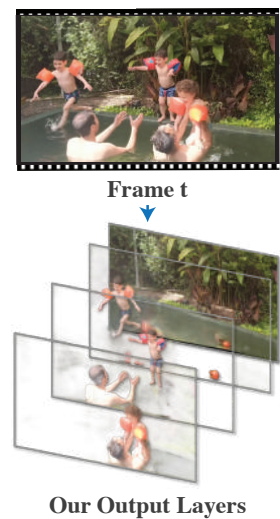


Fig. 1. **Making all children jump into the pool together — in post-processing!** In the original video (left, top) each child is jumping into the pool at a different time. In our computationally retimed video (left, bottom), the jumps of children I and III are time-aligned with that of child II, such that they all jump together into the pool (notice that child II remains unchanged in the input and output frames). In this paper, we present a method to produce this and other people retiming effects in natural, ordinary videos. Our method is based on a novel deep neural network that learns a layered decomposition of the input video (right). Our model not only disentangles the motions of people in different layers, but can also capture the various scene elements that are *correlated* with those people (e.g., water splashes as the children hit the water, shadows, reflections). When people are retimed, those related elements get automatically retimed with them, which allows us to create realistic and faithful re-renderings of the video for a variety of retiming effects. The full input and retimed sequences are available in the supplementary video.

We present a method for retiming people in an ordinary, natural video — manipulating and editing the time in which different motions of individuals in the video occur. We can temporally align different motions, change the speed of certain actions (speeding up/slowing down, or entirely “freezing” people), or “erase” selected people from the video altogether. We achieve these effects computationally via a dedicated learning-based layered video representation, where each frame in the video is decomposed into separate RGBA layers, representing the appearance of different people in the video. A key property

of our model is that it not only disentangles the direct motions of each person in the input video, but also correlates each person *automatically* with the scene changes they generate—e.g., shadows, reflections, and motion of loose clothing. The layers can be individually retimed and recombined into a new video, allowing us to achieve realistic, high-quality renderings of retiming effects for real-world videos depicting complex actions and involving multiple individuals, including dancing, trampoline jumping, or group running.

CCS Concepts: • **Computing methodologies** → **Computational photography**; **Neural networks**.

Additional Key Words and Phrases: retiming, neural rendering, video layer decomposition

ACM Reference Format:

Erika Lu, Forrester Cole, Tali Dekel, Weidi Xie, Andrew Zisserman, David Salesin, William T. Freeman, and Michael Rubinstein. 2020. Layered Neural

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
© 2020 Copyright held by the owner/author(s).
0730-0301/2020/12-ART256
<https://doi.org/10.1145/3414685.3417760>

Rendering for Retiming People in Video. *ACM Trans. Graph.* 39, 6, Article 256 (December 2020), 14 pages. <https://doi.org/10.1145/3414685.3417760>

1 INTRODUCTION

By manipulating the timing of people’s movements, we can achieve a variety of effects that can change our perception of an event recorded in a video. In films, altering time by speeding up, slowing down, or synchronizing people’s motions is often used for dramatizing or de-emphasizing certain movements or events. For example, by freezing the motions of some people in an action-packed video while allowing others to move, we can focus the viewer’s attention on specific people of interest. In this paper, we aim to achieve such effects computationally by retiming people in everyday videos.

The input to our method is an ordinary natural video with multiple people moving, and the output is a realistic re-rendering of the video where the timing of people’s movements is modified. Our method supports various retiming effects including aligning motions of different people, changing the speed of certain actions (e.g., speeding up/slowing down, or entirely “freezing” people). In addition, our method can also “erase” selected people from the video. All these effects are achieved via a novel deep neural network-based model that learns a *layered decomposition of the input video*, which is the pillar of our method. Note that in this paper we focus solely on temporal warping. That is, each person’s pose in our output exists in some frame in the input; we do not generate new, unseen poses or viewpoints.

Motion retiming has been studied so far mostly in the context of character animation, for editing a character’s motion to match a desired duration or target velocity at a given time (e.g., [McCann et al. 2006; Yoo et al. 2015]). In this paper, we take retiming to the realm of *natural real videos*. In the character animation domain, the main challenge is to retime the motion of a set of joints, with the spatiotemporal correlations that exist between them. Analogously, manipulating the timing of people in video not only requires modifying people’s own motions, but also the various elements in the scene that they “cause” and are *correlated* with them — shadows, reflections, the flowing dress of a dancer, or splashing water (Fig. 1). When we retime people, we need to make sure that all those correlative events in the scene follow properly and respect the timing changes. Furthermore, unlike character animation, we do not have any ground truth 3D model of the scene over time; hence, rendering photorealistic, high-quality retiming effects in video is much more challenging. More specifically, retiming motions in videos can often result in new *occlusions* and *disocclusions* in the scene. Rendering the scene content in disoccluded regions and maintaining correct depth ordering between subjects are essential for achieving a realistic effect. Finally, as in any video synthesis task, achieving a temporally coherent result is challenging — small errors such as subtle misalignment between frames immediately show up as visual artifacts when the frames are viewed as a video.

The core of our technique, which we dub *layered neural rendering*, is a novel deep neural-network-based model that is optimized per-video to decompose every frame into a set of layers, each consisting of an RGB color image and an opacity matte α (referred to altogether as “RGBA”). We design and train our model such that each RGBA layer over time is associated with specific people in the video

(either a single person, or a group of people predefined by the user). Crucially, our method does not require dynamic scene elements such as shadows, water splashes, and trampoline deformations to be manually annotated or explicitly represented; rather, only a rough parameterization of the *people* is required, which can be obtained using existing tools with minor manual cleanup for challenging cases. The model then *automatically* learns to group people with their correlated scene changes. With the estimated layers, the original frames of the video can be reconstructed using standard back-to-front compositing. Importantly, retiming effects can be produced by simple operations on layers (removing, copying, or interpolating specific layers) without additional training or processing.

Our model draws inspiration from recent advances in neural rendering [Thies et al. 2019], and combines classical elements from graphics rendering with deep learning. In particular, we leverage human-specific models and represent each person in the video with a single deep-texture map that is used to render the person in each frame. Having a unified representation of a person over time allows us to produce temporally coherent results.

We demonstrate realistic, high-quality renderings of retiming effects for real-world videos with complex motions, including people dancing, jumping, and running. We also provide insights into why and how the model works through a variety of synthetic experiments.

2 RELATED WORK

Video retiming. Our technique applies time warps (either designed manually, or produced algorithmically) to people in the video, and re-renders the video to match the desired retiming. As such it is related to a large body of work in computer vision and graphics that performs temporal remapping of videos for a variety of tasks. For example, [Bennett and McMillan 2007] sample the frames of an input video non-uniformly to produce computational time-lapse videos with desired objectives, such as minimizing or maximizing the resemblance of consecutive frames. [Zhou et al. 2014] use motion-based saliency to nonlinearly retime a video such that more “important” events in it occupy more time. [Davis and Agrawala 2018] retime a video such that the motions (visual rhythm) in the time-warped video match the beat of a target music.

Other important tasks related to video retiming are video summarization (e.g., [Lan et al. 2018]) and fast-forwarding [Joshi et al. 2015; Poleg et al. 2015; Silva et al. 2018], where frames are sampled from an input video to produce shorter summaries or videos with reduced camera motion or shake; or interactive manipulation of objects in a video (e.g., using tracked 2D object motion for posing different objects from a video into a still frame that never actually occurred [Goldman et al. 2008].)

Most of these papers retime *entire video frames*, by dropping or sampling frames. In contrast, we focus on people, and people’s motions, and our effect is applied at the person/layer level. While many methods exist for processing the video in the sub-frame or patch level — both for retiming (e.g., [Goldman et al. 2008; Pritch et al. 2008]) and for various other video manipulation tasks, such as object removal, infinite video looping, etc. (e.g., [Agarwala et al. 2005; Barnes et al. 2010; Newson et al. 2014; Wexler et al. 2007]) —

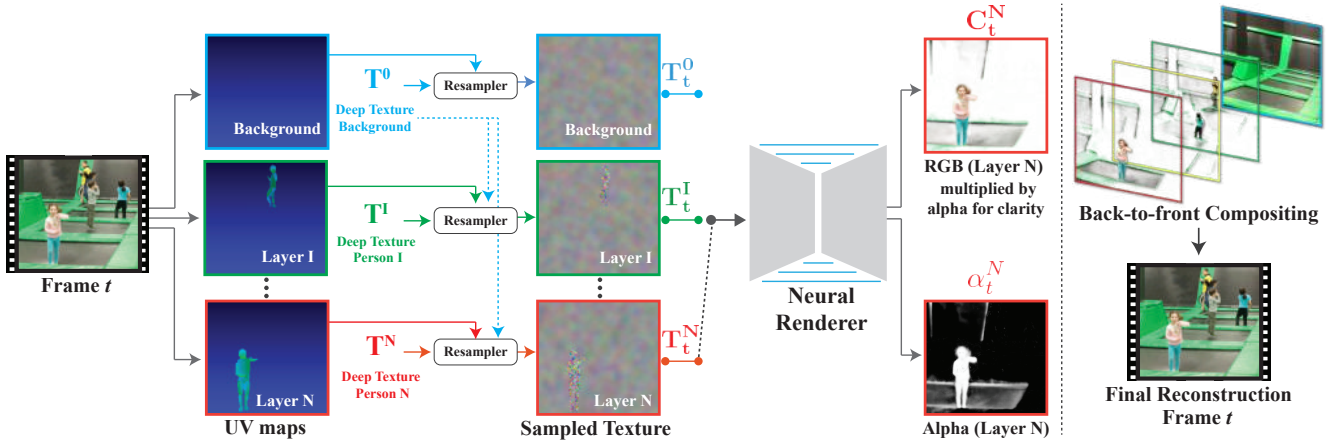


Fig. 2. **Layered Neural Rendering.** Our neural rendering model decomposes each frame of the video into a set of color (C_t^i) and opacity (α_t^i) layers. Each layer is associated with specific people in the video (either a single person, or a group of people predefined by the user). The layers are computed in separate forward passes by feeding to the neural renderer a deep-texture map that corresponds to a single layer. In particular, we represent each person in the video with a single deep-texture map T^i , and the scene background is represented with a deep-texture map T^0 . Given pre-computed UV maps, those deep-texture maps are resampled and composited to form the input to the neural renderer. The set of estimated layers can then be composited in a back-to-front fashion to reconstruct the original frames. Retiming effects can be achieved via simple operations on the layers.

none of these works can handle automatically correlative motions and changes in the video such as shadows and reflections.

Manipulating human poses in video. In the context of manipulating people’s motions, several recent methods have been proposed for transferring motion between two people captured in different videos [Aberman et al. 2019; Chan et al. 2019; Lee et al. 2020; Zhou et al. 2019] or transferring motion from a low-dimensional signal [Gafni et al. 2020]. However, there are two main distinctions between these methods and ours: (i) motion transfer methods are focused on “puppetry” – generating people in new unseen poses – whereas our method is tailored solely for time warping (we do not generate new poses); and (ii) existing motion transfer methods consider only a single person in each video – they do not model occlusions/disocclusions between multiple people – and with the exception of [Gafni et al. 2020], they also do not capture the correlations between people and their environment, elements that are crucial for producing realistic retiming effects.

Image and video matting. Our work is also related to image and video matting as we decompose each input video frame into a set of RGBA layers. Nevertheless, both traditional matting techniques (e.g., [Bai et al. 2009; Chuang et al. 2002; Li et al. 2005; Wang et al. 2005]) as well as more recent learning-based matting methods [Hou and Liu 2019; Xu et al. 2017] do not capture correlated effects such as shadows and reflections. Moreover, matting methods typically rely on accurate trimaps that are often given manually by the user, whereas we automatically generate rougher trimaps from detected keypoints. Even though our trimaps are not as accurate, we are still able to learn complex effects such as loose clothing, even when such regions are not included in the initial, estimated trimaps. Finally, existing matting methods cannot handle well entirely semi-transparent objects like reflections (in matting trimaps, pixels that

are labeled “opaque” are fully assigned to the foreground). We show comparisons with image matting in Section 6.

Layered representations for video. Decomposing videos into layers is a classical problem in computer vision originally proposed in [Wang and Adelson 1994], which has inspired a large body of work (e.g. [Fradet et al. 2008; Jovic and Frey 2001; Kumar et al. 2008; Nandoriya et al. 2017; Xue et al. 2015; Zitnick et al. 2004]). Recently, several works have leveraged the power of deep neural networks to decompose videos into layers. For example, [Zhou et al. 2018] and [Srinivasan et al. 2019] train a network on input stereo pairs of images to predict layers that correspond to physically accurate depth, with the goal of view synthesis for static scenes. Our layers are not meant to be strictly accurate in terms of depth, as we do not aim to generate scenes from novel viewpoints, but rather composite scene elements from different points of a dynamic video.

Alayrac et al. [2019a; 2019b] train a model to decompose synthetically blended real videos and then apply it to natural videos to remove reflections, shadows, and smoke. Their method is general rather than person-specific like ours, and they use audio as a cue for separating layers, rather than people geometry as we do.

Finally, Gandelsman et al. [2019] extend the idea of “Deep Image Prior” (DIP) networks [Ulyanov et al. 2018] for decomposing a single image or a video into two layers. Their key observation is that the structure of a CNN provides a prior that drives each layer towards a natural image. As above, they have limited control over their output layers and rely entirely on CNN properties to produce meaningful layers. Our method also leverages this property of CNNs (see Section 5), yet provides control over the decomposition, which is required for retiming effects. Furthermore, their method is designed for two-layer decomposition, and generalizing it to N layers requires $N \times$ learnable parameters. In contrast, our method can produce an arbitrary number of layers with the same, fixed number of

network parameters. (We show a comparison with Double-DIP in the supplementary material.)

Neural rendering. Neural networks have recently begun to be used as a final rendering layer for 3D scenes [Kim et al. 2018; Liu et al. 2019; Martin-Brualla et al. 2018; Meshry et al. 2019; Sitzmann et al. 2019; Thies et al. 2019]. A neural network can be trained to bridge the gap between incomplete or inaccurate input geometry and a photorealistic result. Examples of incomplete input are point clouds [Meshry et al. 2019], voxels [Sitzmann et al. 2019], partial 3D scans or 3D models of humans [Kim et al. 2018; Liu et al. 2019; Martin-Brualla et al. 2018], and textured proxy geometry [Thies et al. 2019]. Our work adopts the textured proxy approach – creating geometric proxies for each person in the scene and texturing them with a deep texture map [Thies et al. 2019]. Importantly, we combine this approach in a neural rendering model that outputs a layered decomposition of each frame of the video. This is in contrast to existing neural rendering methods, which output the final reconstruction directly. Our layered decomposition is the key to producing realistic high-quality retiming effects, as demonstrated in Fig. 11 and discussed in detail in Section 6.4.

3 OVERVIEW

A key challenge in generating realistic retiming effects in videos is ensuring that when people’s motions are retimed, their related effects in the scene follow with them. For example, if we freeze a person, their shadow should freeze as well.

We achieve this via a learning-based approach, which, at its core, decomposes each frame of the input video into a *human-specific set of layers*. That is, each layer is associated with specific people in the video—either a single person, or a group of people (Section 4.1). A key feature in our approach is that only the people in the video are modeled explicitly, while the rest of the scene elements correlated with each person are *inferred automatically* by our layered neural rendering pipeline.

Our pipeline is illustrated in Fig. 2. Our model is trained per-video, observing only the input video (without any additional examples), and is trained in a self-supervised manner, by learning to reconstruct the original video frames by predicting the layers (Section 4.2). We use off-the-shelf methods (AlphaPose [Fang et al. 2017], DensePose [Güler et al. 2018]) in combination with our own techniques to represent each person at each frame (Section 4.5). This representation is then passed on to the neural renderer and is used for seeding the layer assignment (person per layer). While the input to the neural renderer includes only the people (and a static background), the renderer’s task is to generate layers that reconstruct the *full* input video. Thus, it must assign all the remaining, time-varying scene elements (e.g., shadows, reflections) into the appropriate peoples’ layers.

The neural renderer succeeds in this task because the network design and training procedure (Section 4.3) encourages scene elements that are correlated with a layer to be captured faster (reconstructed earlier in training) than non-correlated elements. This property is related to the “Deep Image Prior” (DIP) work [Ulyanov et al. 2018], which showed that overfitting a network to a natural image required fewer gradient descent steps than fitting a noise image. We further

explore how correlated effects are captured in the correct layers through synthetic experiments in Section 5.

4 METHOD

4.1 A Layered Video Representation

Given an input video V , our goal is to decompose each frame $I_t \in V$ into a set of RGBA (color + opacity) layers:

$$\mathcal{L}_t = \{L_t^i\}_{i=1}^N = \{C_t^i, \alpha_t^i\}_{i=1}^N, \quad (1)$$

where C_t^i is a color image and α_t^i is an opacity map (matte). The i^{th} layer for all frames L_t^i is associated with person i in the video. We add an additional background layer L_t^0 , not associated with any person, that learns the background color.

Given this layered representation and a back-to-front ordering for the layers, denoted by o_t , each frame of the video can be rendered using the standard “over” operator [Porter and Duff 1984]. We denote this operation by:

$$\hat{I}_t = \text{Comp}(\mathcal{L}_t, o_t) \quad (2)$$

We assume that the compositing order o_t is known, yet time varying, i.e., the depth ordering between people may change throughout the video (see Section 4.5).

A key property of this representation is that retiming effects can be achieved by simple operations on individual layers. For example, removing person i from frame t can be done simply by removing the i^{th} layer from the composition, i.e., by substituting in $\mathcal{L}_t \setminus L_t^i$ into Eq. 2. Similarly, generating a video where person i is frozen at a time t_0 is achieved by copying $L_{t_0}^i$ over L_t^i for all frames. We expand on other operations in Section 6.

4.2 Layered Neural Rendering

Decomposing a real-world video into a desired set of layers is a difficult, under-constrained problem – there are numerous possible decompositions that can provide an accurate reconstruction of the original frame I_t . For example, a single visible layer that contains the entire frame can perfectly reconstruct the video. We therefore constrain our layered neural renderer to steer the solution towards the desired person-specific decomposition. We do so in several ways including incorporating human specific representations, tailoring the input to the network, and applying dedicated losses and training regimes.

As mentioned in Section. 3, the layers are predicted in *separate feed-forward passes* through the neural renderer, where the input for each pass represents only one person (or a pre-defined group of people). In doing so, we control the association of people with the output layers. More specifically, we construct the input to the renderer as follows:

Person representation. We parameterize each person in the video with a single human texture atlas T^i and a per-frame UV-coordinate map UV_t^i , which maps each pixel in the human region in frame I_t to the texture atlas. We use the parameterization of the SMPL model [Loper et al. 2015] that can be estimated from an input image using existing methods (e.g., DensePose [Güler et al. 2018]). This provides a unified parameterization of the person over time and a convenient model for appearance and geometry. We follow a



Fig. 3. **Example layer decompositions by our model.** For each example, we show (a) an original frame from the video, (b) the predicted background layer, and (c) two representative (RGBA) layers associated with people (in *Ballroom* and *Splash*, third and fourth rows, respectively, the grouping of people into layers was specified manually; in the other examples each person is assigned to a separate layer). Our model successfully disentangles the people into different layers along with the visual changes they cause on the environment, such as trampoline deformations, shadows, reflections, and water splashes. (Artifacts at the bottom of the background layer in the pool scene (bottom row) are due to the fact that those regions are never visible in the input video.)

similar approach to Thies *et al.* [2019] and replace the classic RGB texture map with a learnable, high-dimensional texture map, which can encode more powerful and richer appearance information. The deep texture maps are then decoded into RGB values using a neural rendering network. To represent person i at time t , we sample its deep texture map T^i using UV_t^i , obtaining T_t^i .

Background representation. The background is represented with a single texture map T^0 for the entire video. Sampling from the background is performed according to a UV map UV_t^0 . For a static camera, UV_t^0 is an identical xy coordinate grid for all frames. If homography transformations estimated from camera tracking are available (see Section 4.6), UV_t^0 is the result of transforming an xy coordinate grid by the homography for frame t .

Input-Output. To construct the input, we place the background's UV map behind each person's UV map to provide background context for the renderer (see Fig. 2). That is, the input for layer i at time t is the sampled deep texture map T_t^i , which consists of person i 's sampled texture placed over the sampled background texture. Each

of the re-sampled texture layers $\{T_t^i\}_{i=1}^N$ as well as the background re-sampled texture is fed to the renderer separately. The output is $L_t^i = \{C_t^i, \alpha_t^i\}$, the time-varying color image and opacity map for that layer, respectively.

The neural renderer is trained to reconstruct the original frames from the predicted layers, using Eq. 2. The depth ordering of the people layers, which may vary over time, is provided by the user. To accurately reconstruct the *entire* original video, the neural renderer must reconstruct any space-time scene element that is not represented by the input UV maps, such as shadows, reflections, loose clothing or hair. Crucially, disentangling people into separate inputs is key for guiding the neural renderer to assign those scene elements into the layers with which they are most strongly correlated in space and/or time (e.g. attached to a person, moves like the person, etc.). We expand on this in greater detail in Section 5.

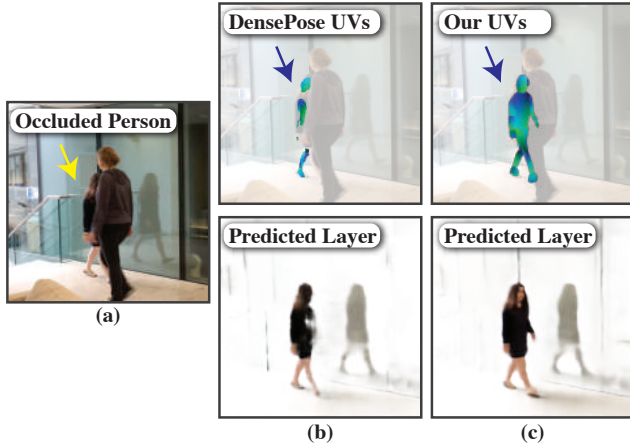


Fig. 4. **Full-body UVs.** Existing tools such as DensePose [Güler et al. 2018] provide UV coordinates only for fully visible regions. In contrast, our method is able to produce full-body UVs for occluded people by detecting and interpolating keypoints, then passing the occluded person’s full skeleton to a trained network that converts it to a UV map. This allows us to create editing effects that require disoccluding people.

4.3 Training

We now turn to the task of learning the optimal parameters θ of the neural renderer, and the set of latent textures $\{T^i\}_{i=0}^N$ by optimizing the learned decomposition for each frame.

One necessary property of the learned decomposition is that it will allow us to accurately reconstruct the original video. Formally,

$$E_{\text{recon}} = \frac{1}{K} \sum_t \|I_t - \text{Comp}(\mathcal{L}_t, o_t)\|_1, \quad (3)$$

where \mathcal{L}_t are the output layers for frame t , o_t is the compositing order, and K is the total number of frames.

The reconstruction loss alone is not enough to make the optimization converge from a random initialization, so we bootstrap training by encouraging the learned alpha maps α_t^i to match the people segments that are associated with layer i . To do so, we apply the following loss:

$$E_{\text{mask}} = \frac{1}{K} \frac{1}{N} \sum_t \sum_i D(m_t^i, \alpha_t^i), \quad (4)$$

where m_t^i is a trimap derived from the UV maps UV_t^i (see Fig. 10), and $D(\cdot)$ is a distance measure. Masks are trimaps with values in $[0, 0.5, 1]$, where the uncertain area is produced by morphological dilation of the binary UV mask. For a trimap m , let b_0 be the binary mask of the pixels where $m = 0$, with b_1 defined likewise. Losses in the positive and negative regions are balanced while the uncertain area is ignored. The distance measure is:

$$D(m, \alpha) = \frac{\|b_1 \odot (1.0 - \alpha)\|_1}{2 \|b_1\|_1} + \frac{\|b_0 \odot \alpha\|_1}{2 \|b_0\|_1} \quad (5)$$

where \odot is the element-wise (Hadamard) product.

Since the UV mask does not include information from correlated effects such as shadows and reflections, E_{mask} is only used to bootstrap the model and is turned off as optimization progresses.

We further apply a regularization loss to the opacities α_t^i to encourage them to be spatially sparse. This loss is defined as a mix of L_1 and an approximate- L_0 :

$$E_{\text{reg}} = \frac{1}{K} \frac{1}{N} \sum_t \sum_i \gamma \|\alpha_t^i\|_1 + \Phi_0(\alpha_t^i) \quad (6)$$

where $\Phi_0(x) = 2 \cdot \text{Sigmoid}(5x) - 1$ smoothly penalizes non-zero values of the alpha map, and γ controls the relative weight between the terms.

Our total loss is then given by:

$$E_{\text{total}} = E_{\text{recon}} + \gamma_m E_{\text{mask}} + \beta E_{\text{reg}}, \quad (7)$$

where γ_m and β control the relative weights of the terms.

As usual with non-linear optimization, the results are sensitive to the relative weights of the error terms. To produce the results in this paper, we vary the loss weights based on the current epoch e and the “bootstrap” epoch e_b , which is the first epoch where $E_{\text{mask}} < 0.02$:

$$\gamma_m = \begin{cases} 50 & e \leq e_b \\ 5 & e_b < e \leq 2e_b \\ 0 & \text{otherwise} \end{cases} \quad \gamma = \begin{cases} 2 & e \leq 200 \\ 0 & \text{otherwise} \end{cases} \quad \beta = 0.005 \quad (8)$$

This schedule puts a heavy initial loss on the masking term to force the optimization towards a plausible solution, then relaxes it to allow the optimization to introduce effects not present in the masks.

4.4 High-Resolution Refinement and Detail Transfer

We take a multi-scale approach and first train our model using Eq. 7 on a downsampled version of the original video. We then upsample the result to the original resolution using a separate lightweight refinement network, which consists of several residual blocks (see Appendix for details) operating on each RGBA layer separately.

We can avoid the additional expense of training with perceptual and adversarial losses by directly transferring high-resolution details from the original video in a post-processing step. The residual between the neural renderer output and the video defines the detail to transfer, and the amount of the residual to transfer to each layer is determined by the transmittance map τ_t^i :

$$\tau_t^i = 1.0 - \text{Comp}_\alpha(\mathcal{L}_t \setminus \{L_t^j \mid j < i\}, o_t \setminus \{j \mid j < i\}) \quad (9)$$

where Comp_α denotes the alpha channel of the composite produced by the neural renderer. The final layer colors are:

$$C_t^i = \text{Cnr}_t^i \cdot \tau_t^i (I_t - \text{Comp}(\mathcal{L}_t, o_t)) \quad (10)$$

where Cnr is the color produced by the neural renderer. See supplementary material (SM) for visualizations. Given this transfer, the upsampling network needs only to refine the predicted alpha mattes and produce reasonable colors in occluded regions, where ground-truth high-frequency details are not available. For pixels with no occlusion and a predicted alpha value of 1, this high-frequency transfer amounts to using the original RGB pixel.

4.5 Keypoints-to-UVs

Estimating people UV maps from images can be done via off-the-shelf methods (e.g., DensePose [Güler et al. 2018]). However, using such methods directly has two main drawbacks:

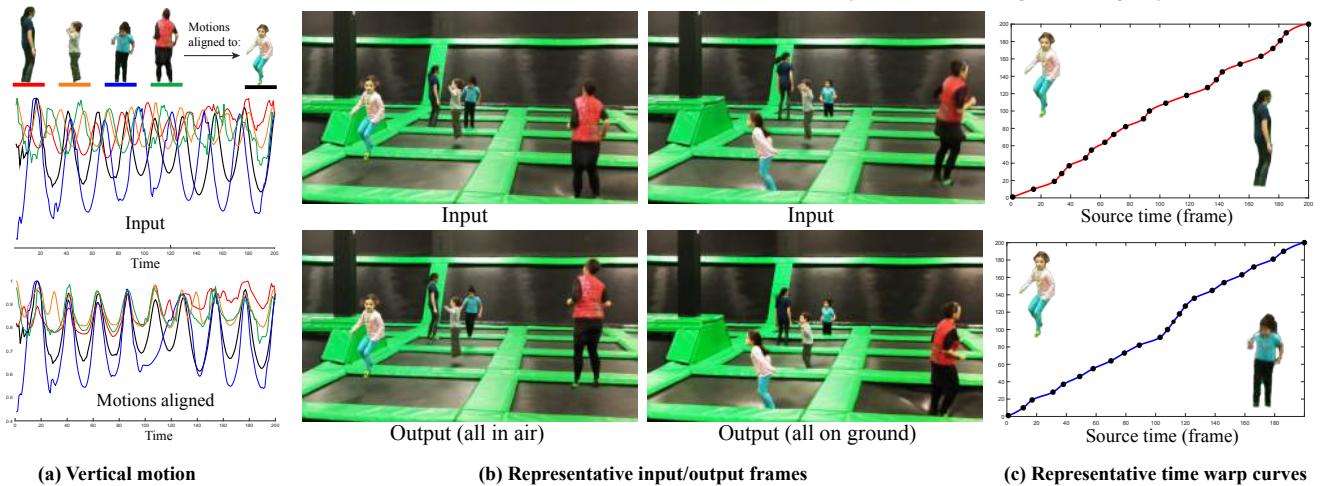


Fig. 5. **Automatic motion alignment.** People jumping on trampolines are motion-aligned to the little girl in the front left, such that they all jump in sync. (a) Top: motion signals for the different people in the video, taken as the normalized y -coordinate of the center of mass of each person over time, as computed from each person's tracked keypoints. Bottom: time-warped motion signals, after aligning each one to the motions of the little girl (black line in the plots) using Correlation Optimized Warping [Tomasi et al. 2004] (a variation of Dynamic Time Warping). (b) Each column shows corresponding frames from the input video (top), and from the output, retimed video (bottom). (c) Two estimated time warps, aligning two of the jumpers with the jumps of the front little girl.

- (1) *Disoccluded people:* Video retiming effects often result in disocclusions of people who were partially or even fully occluded in the original frames. Because we model the appearance of a person with a single texture map that is learned jointly for the entire video, we can render disoccluded content as long as we can correctly sample from it. Thus, we want to ensure that all UV maps represent the full body of each individual, even in the presence of occlusion. Existing methods, however, provide an estimate only for the visible human regions (Fig. 4(a)).
- (2) *Robustness:* Existing methods tend to suffer from erratic errors such as missing body parts, especially in the presence of motion blur, noise or occlusion (see SM for examples). Due to their dense nature, temporally filtering or interpolating UV maps in order to remove such errors is challenging.

To overcome these limitation, we train a neural network that predicts clean full-body UV maps from human keypoints. Keypoint estimators such as AlphaPose [Fang et al. 2017] tend to be more robust, and unlike UV maps, can easily be interpolated and manually corrected. To train our keypoint-to-UV model, we use the *Let's Dance Dataset* [Castro et al. 2018], curated to contain only single-person video frames, and our own filmed video of approximately 10 minutes of a single person doing a variety of poses. We then generate approximately 20K keypoint-UV training examples in which people are fully visible, by running AlphaPose and DensePose on the original frames. We follow a similar training regimen as in [Güler et al. 2018]; see Appendix for further details.

At test time, we can estimate full-body UV maps even in the presence of significant occlusions: we first estimate keypoints using AlphaPose, track the keypoints using PoseFlow [Xiu et al. 2018], and linearly interpolate occluded keypoints. The trained keypoint-to-UV network then processes these keypoints to generate complete UV maps (Fig. 4(b)). In challenging cases, current state-of-the-art

keypoint detectors and trackers can still fail due to motion blur, occlusions, or noise. For the results in this paper, we manually fixed these errors using a rotoscoping tool for between zero and 10% of input frames, taking a maximum of 30 minutes for particularly troublesome videos. See SM videos for examples of keypoint edits.

4.6 Camera Tracking

When the input video contains a moving camera (as is the case with all our video examples), we first estimate the camera motion using a feature-based tracking algorithm similar to the one described in [Grundmann et al. 2011]. We model the camera motion at each frame using a homography transformation, which we estimate robustly from matched ORB features [Rublee et al. 2011] between frames. When stabilizing small camera motions or natural hand shake, we compute the homographies between each frame and a single reference frame (which works better than tracking the camera over time), then use them to stabilize the input video. When the video contains large camera motion or substantial panning, we estimate homographies over time between consecutive frames, use them to register all the frames with respect to a common coordinate system, then apply this coordinate system transformation to the background UV map to preserve the original camera motion. Retiming a layer from frame t to \bar{t} is achieved by transforming the layer to the common coordinate system using the transformation at t , then applying the inverse transformation at \bar{t} .

5 WHY IT WORKS

As discussed in Section 4.3, our training ensures that specific people are assigned into specific layers (by bootstrapping the predicted alpha mattes with rough people trimaps), whereas effects outside the human trimap region, such as shadows and reflections, are learned incrementally as the network trains. But how does the



Fig. 6. **Synthetic examples for testing effects of motion correlation and spatial proximity.** In video 1, one square moves smoothly with the person while the other is randomly placed. The smoothly moving square is reconstructed earlier in training than the random square. In video 2, two squares have the same motion but one is closer to the person. The closer square is reconstructed earlier in training.

model manage to associate such scene elements with the correct layers?

Single person and their correlated effects. The reason for this successful assignment stems from properties of CNNs. Specifically, space-time scene elements that are most strongly correlated with the input are reconstructed earlier in training (i.e., require smaller changes to the network weights) than non-correlated elements. To illustrate this, consider the synthetic sequences shown in Fig. 6 (the full videos are available in the SM). A person is superimposed on a fixed background along with two green squares. The first video examines correlations in motion and the second examines correlations in space (proximity). In the first video, one square moves consistently left-to-right behind the person, while the other is randomized around the person. The network first learns to reconstruct the static background and person, then the correlated square, and finally the random square. In the second video, two squares have the same motion, but one is closer to the person. The closer square is reconstructed earlier in training than the farther square.

This behavior is tightly related to the “Deep Image Prior” (DIP) principle [Ulyanov et al. 2018]. An image is represented by the parameters of a CNN by “overfitting” it with gradient descent until the network output matches the image. Importantly, it was shown that fitting a natural image (which contains repetitive patches) takes

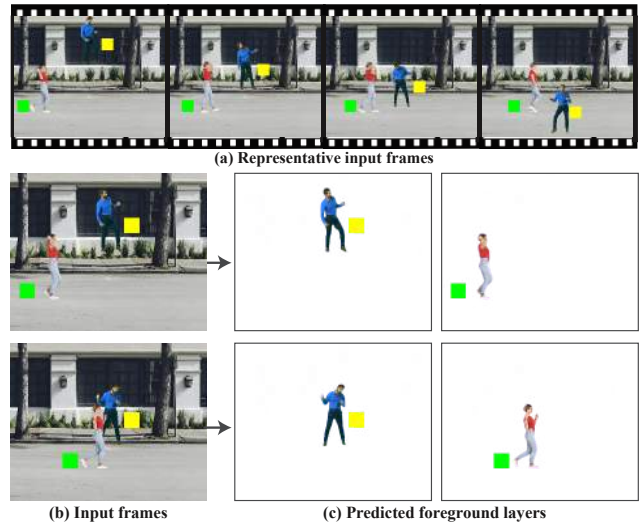


Fig. 7. **A synthetic example demonstrating assignment into layers.** (a-b) A woman (in red) is translating from left to right together with a green square next to her. Similarly, a man (in blue) is oscillating up and down with a yellow square next to him. (c) With gradient descent iterations, each person layer “grabs” more quickly the square that is closest and moves most similarly to it (see Fig. 6). The squares are consistently assigned into each person’s layers throughout the video, even in the presence of occlusions between the people and squares.

fewer iterations than fitting a random noise image (or a natural image corrupted by noise). In our case, effects that are closely correlated with a person in space, time, or both, are learned more quickly than other effects. Our synthetic examples and results on natural videos support this observation.

Assigning the correct effects to multiple people. When there are multiple people present in the video, each person layer needs to “grab” the scene elements that are correlated with it. This assignment is achieved via our training regime and the construction of the input to the network. The renderer sees one layer’s UV coordinates at a time, so must explain each missing scene element using only the information from a single person (or preselected group of people). Because correlated effects are learned faster than uncorrelated effects, each person layer will “grab” the effects that are most correlated with it before the other layers do. Once an effect is reconstructed in a layer, the reconstruction loss for that effect is minimized and no gradient encourages the effect to be learned in other layers. In practice, this means that loose clothing, shadows, reflections, or other effects that are most strongly correlated with the person creating them, are likely to be decomposed into the correct layer.

Fig. 7 shows such a synthetic, two-person example. The second green square from the synthetic sequence in Fig. 6 is replaced by a second person (a man in blue), who has a yellow correlated square. The motions of the two people are uncorrelated: a man (in blue) oscillates vertically while a woman (in red) translates horizontally. We bootstrap each layer with its respective person’s trimap and train the network to reconstruct the video. The network learns to reconstruct the green square in the woman’s layer, and the yellow



Fig. 8. **Retiming and editing results** for person removal (top) and motion freeze (bottom). For each example (pairs of rows) we show sample frames from the original video (top row) and their corresponding retimed/edited frames (bottom row).

square in the man’s layer. Further, the yellow square (and the man) are fully reconstructed, even for frames where they were occluded by the woman or the green square.

Note that if instead of the layered input, the network sees the entire scene at one time (similar to [Thies et al. 2019]), the network learns to explain the missing scene elements using *combinations* of people, not individuals. Artifacts then appear when individuals are retimed to form a new combinations not seen during training (Fig. 11).

6 RESULTS

We tested our method on a number of real-world videos, most of which are captured by hand-held cellphone cameras (either captured by us or taken from the Web). All the videos depict multiple people moving simultaneously and span a wide range of human actions (e.g., dancing, jumping, running) in complex natural environments. Representative frames from these videos are shown in Figs. 1, 3

and 8, and the full input and output sequences are available in the supplementary material.

We implemented our networks in PyTorch and used the Adam optimizer [Kingma and Ba 2014] with a learning rate of $1e - 3$. Depending on video length and number of predicted layers, total training time took between 5 and 12 hours on 2 NVIDIA Tesla P100 GPUs. See Appendix A.1 for further implementation details.

6.1 Layer Decomposition

Several of our layer decompositions are visualized in Fig. 3. For *Ballroom* and *Splash*, we group certain people into one layer; for all other videos, each person is assigned with their own layer. For all the videos, our model successfully disentangles the people into the layers. The layers capture fine details such as loose hair and clothing (e.g., foreground dancer’s white dress in *Ballroom*), or objects attached to the people (e.g., children’s floaties in *Splash*). This is in spite of initializing our model with a rough people UV map without explicitly representing these elements. This ability of our

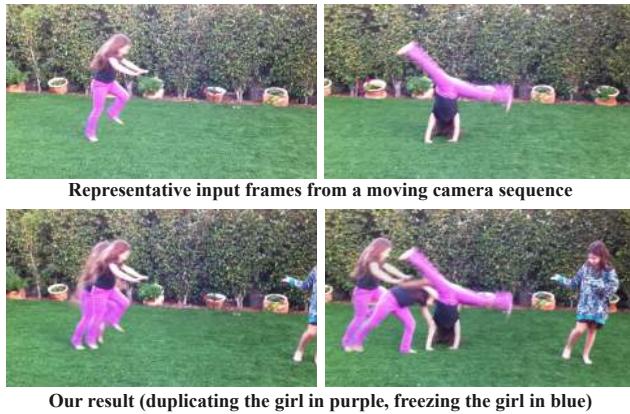


Fig. 9. **Cartwheel.** In the original video, the camera is panning. In our edited result, we preserve the original camera motion, while duplicating the girl who performs the cartwheel and freezing the girl in blue so that she remains in the frame rather than moving offscreen.

model to accurately segment the people regions is also illustrated more closely in Fig. 10.

Furthermore, the predicted layers successfully correlated the people with other nontrivial visual time-varying scene elements that are related to them—for example, shadows cast on the floor by the different dancers (*Ballroom*), complicated reflections of two people crossing each other (*Reflections*), surface deformation (*Trampoline*), or water splashes (*Splash*) caused by people’s motion.

6.2 Retiming and Editing Results

With the predicted layers in hand, we produced a variety of retiming and editing effects via simple operations on the layers. We show several such retiming results in Fig. 1, Fig. 5, and Fig. 8 (full videos are available in the SM).

In Fig. 1 and Fig. 5, we have multiple people performing a similar action, but their motions are not in sync. In Fig. 1, the children jump into the pool one after the other, whereas in Fig. 5 the periodic motions of the people jumping are independent. In both examples, we retime the people to align their motions. For *Splash*, several alignment points defined manually were sufficient to align the children’s jumps. In *Trampoline*, because of the periodic nature of the motion, the alignment is performed automatically using Correlation Optimized Warping [Tomasi et al. 2004] (a variation of Dynamic Time Warping). Notice how all the time-varying scene elements caused by the people—water splashes as they hit the water, trampoline deformations as they bounce on it—follow automatically with them in our retimed results.

We can also use our method to “freeze” people at a certain point in time, while letting other people move as in the original video. This allows viewers to focus their attention on the moving people while ignoring the rest of the motions in the original video. In *Ballroom* (Fig. 8), we freeze the dancing couple in the back throughout the video, while the couple in front keeps moving. Here too, the shadows and reflections on the floor move realistically with the moving couple, while the shadows of the background couple remain static.

Furthermore, disoccluded regions of the back couple are realistically rendered.

In *Kids Running* (supplementary material), we show how our model can scale to multiple layers (8) to produce complex retiming effects involving many people. We retime the original video, where the kids are crossing the faint finish line on the ground at different times, to produce a ‘photo-finish’ video, where all the kids cross the finish line together. We achieve this result by slowing down the layers of the children that run offscreen. Even though this sequence involves many individuals, our model is able to obtain clean mattes for each child. In addition, occluded people and large occluded regions in the background are realistically inpainted, all while handling significant motion blur that exists in the input video.

As mentioned, in addition to retiming effects, our method can also support easy removal of people in video—a byproduct of the layered representation we use. In *Reflections*, we demonstrate person removal in a video containing two people crossing paths in front of a window. Here our model manages to perform several nontrivial tasks: it completely disoccludes the person walking in the back; it associates each person properly with their reflection and shadow; and it disentangles the two reflections when they overlap (despite the fact that none of these elements are represented explicitly by the model). Consider generating such a result with a traditional video editing pipeline: the reflections would have to be tracked along with the people to perform proper removal; the person in the back would have to be manually inpainted at the point where they are occluded by the person in front. The advantage of our method is that by merely inputting UVs for each person in a separate layer, and turning those layers ‘on and off’, we can achieve the same result with significantly less manual work.

In Fig. 9, we show an example where the original camera motion is preserved in our retimed result (see Section 4.6) – the girl on the left is duplicated with a short time offset between each of her copies, and the girl on the right is frozen, all while the camera is panning as in the original video.

6.3 Comparisons

We compare our method to state-of-the-art learning-based image matting [Hou and Liu 2019] and Double-DIP [Gafni et al. 2020] in Fig. 10. Matting methods are unsuitable for retiming because they fail to capture desired effects outside of the trimap region. To capture such effects using matting would require annotating the trimap to represent these regions. Our method, however, is not strictly bound by the trimap, as we use it to bootstrap our training only in the initial epochs, after which the optimization is allowed to grow the matte to encompass regions beyond the trimap border. Thus, we can capture complex effects *without explicitly representing them in the trimap*, merely by representing the human region coarsely using automatically generated trimaps. While Double-DIP is similarly unbound by the trimap region, it fails to segment the entire bicycle.

6.4 Ablations

We ablate components of our method, specifically the keypoint-to-UV network and layer decomposition.

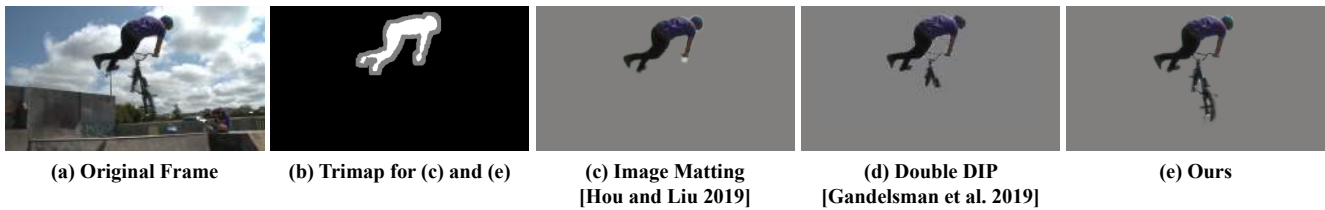


Fig. 10. **Comparison with matting and Double-DIP.** We obtain trimaps from a mask of the person’s UV; the gray pixels represent a dilated region around this mask for which the matching loss is ignored. Matting is insufficient for our retiming task because it fails to capture correlated effects that are not included in the input trimap (e.g. the bicycle). Unlike matting, our final model prediction is able to move beyond the bounds provided by the coarse trimap to produce a more accurate segmentation mask that includes correlated regions. While Double-DIP is not limited by trimap bounds, it fails to segment the entire bicycle.

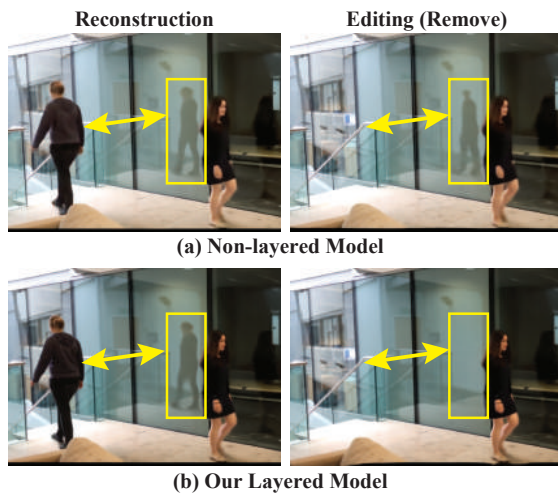


Fig. 11. **Layered vs. non-layered neural rendering.** Removing the layer decomposition component of our model results in a lack of generalization ability. While both models are able to reconstruct the original video, the non-layered model struggles to disentangle people and their correlations (reflection in right column is not fully removed when the person is removed).

6.4.1 Keypoint-to-UV network vs. DensePose. We compare layer decomposition results from using our keypoint-to-UV network to predict UV maps versus using DensePose outputs directly (Fig. 4). In the case of partial occlusion, our UV maps manage to complete the occluded regions of the person, while the DensePose UVs fail to complete the face and the left arm. In a more extreme case where the person is fully occluded, but their reflection is visible, it is crucial to generate a UV map for that person not only for the purposes of disocclusion, but because any visible effects such as their reflection or shadows will be incorrectly reconstructed in the layer of the visible person (as the fully occluded UV will be indistinguishable from the background-only UV). Thus we disocclude UVs by using keypoints as an intermediate representation that allows for more robust estimation and straightforward interpolation.

6.4.2 Layered vs. Non-Layered Neural Rendering. We compare our method to a non-layered model, i.e., a model that takes as input a single sampled texture map representing all the people in the frame as well as the background, and directly outputs an RGB reconstruction of the frame. This baseline follows the deferred neural rendering approach [Thies et al. 2019].

Fig. 11 shows a reconstruction and editing result produced by the non-layered model. As can be seen, the non-layered model can reconstruct the original frames fairly well despite the missing information and noise in the UV maps. This aligns with the results demonstrated in [Thies et al. 2019]. However, when editing is performed, the non-layered model performs poorly and is unable to generalize to new compositions of people, as evident by the partial reflection artifact. The main reason is twofold: (i) To produce editing effects, the model is required to generalize to new UV compositions of people in configurations that were never seen during training; thus it struggles to produce realistic-looking results based solely on L1 reconstruction loss. Our approach avoids this generalization issue because editing is performed as post-processing on the predicted layers—the same outputs produced during training. (ii) When the input to the model is a composition of all the people in the frame rather than separated UVs, the model can easily reconstruct the original frame without necessarily capturing meaningful correlations since it is not required to disentangle separately moving parts of the scene. This can be seen where the non-layered model struggles to learn the correct relationships between the different people and their reflections.

Another benefit of the layered representation is that we can transfer high-resolution details from the input video to each layer, as described in Section 4.4. Because the non-layered model produces the retimed result directly, there is no opportunity to transfer detail from the input video prior to retiming. Achieving comparable quality with the non-layered model would likely require extensive training time and additional loss functions.

6.5 Limitations

While our system manages to decompose scenes quite well—even surprisingly so—in some scenes the results may not be perfect. Here we point out some of the limitations that we have observed.

For a particularly difficult section of the *Splash* video, part of the large water splash caused by the child who is underwater in *Layer 1* is incorrectly placed in *Layer 2* (Fig. 12). However, this error can be fixed by manually editing the initial trimap masks for *Layer 1* to include the entire splash (top right, Fig. 12). Such editing can be done very quickly when needed, requiring only rough marking on a few of the frames (a single trimap in this example).

Another type of artifact that may occur is when the scene includes time-varying background elements that are not correlated with any of the people of interest in the video. For example, since we assume a static background, the colorful blinking lights in the *Ballroom* video

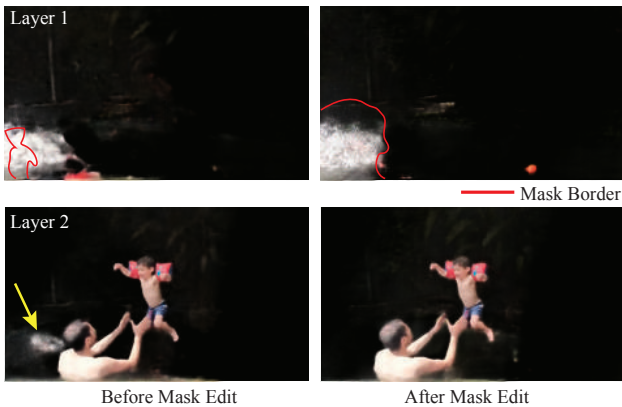


Fig. 12. **User-guided trimaps.** In some particularly difficult segments, the assignment of effects to layers may be incorrect, e.g., some parts of the splash caused by the child who is underwater are incorrectly assigned to Layer 2 (yellow arrow). The user can correct such errors by simple manual editing of the initial trimaps. Here, we expand the initial trimap to include the entire splash region (red outline); this was done for only one layer of one frame, and then duplicated for a segment of the video.

must be included in one of the foreground layers. As can be seen in Fig. 3, the purple light appears in the front couple’s layer, making it impossible to retime them separately. This artifact could be remedied by adding a separate layer to represent the dynamic background elements in addition to the layers for the dancing couples.

7 CONCLUSION

We have presented a system for retiming people in video, and demonstrated various effects, including speeding up or slowing down different individuals’ motions, or removing or freezing the motion of a single individual in the scene. The core of our technique is learned layer decomposition in which each layer represents the *full* appearance of an individual in the video – not just the person itself but also all space-time visual effects correlated with them, including the movement of the individual’s clothing, and even challenging semi-transparent effects such as shadows and reflections.

We believe that our layered neural rendering approach holds great promise for additional types of synthesis techniques, and we plan to also generalize it to other objects besides people, and to expand it to other non-trivial post-processing effects, such as stylized rendering of different video components.

ACKNOWLEDGMENTS

We thank the friends and family that appeared in our videos. The original *Ballroom* video belongs to Desert Classic Dance. This work was funded in part by the EPSRC Programme Grant Seebibyte EP/M013774/1.

REFERENCES

Kfir Aberman, Mingyi Shi, Jing Liao, Dani Lischinski, Baoquan Chen, and Daniel Cohen-Or. 2019. Deep Video-Based Performance Cloning. In *Computer Graphics Forum*, Vol. 38. Wiley Online Library, 219–233.

Aseem Agarwala, Colin Zheng, Chris Pal, Maneesh Agrawala, Michael Cohen, Brian Curless, David Salesin, and Richard Szeliski. 2005. Panoramic Video Textures. In *SIGGRAPH*.

Jean-Baptiste Alayrac, João Carreira, and Andrew Zisserman. 2019a. The Visual Centrifuge: Model-Free Layered Video Representations. In *CVPR*.

Jean-Baptiste Alayrac, Joao Carreira, Relja Arandjelovic, and Andrew Zisserman. 2019b. Controllable Attention for Structured Layered Video Decomposition. In *ICCV*.

Xue Bai, Jue Wang, David Simons, and Guillermo Sapiro. 2009. Video SnapCut: robust video object cutout using localized classifiers. *TOG* (2009).

Connelly Barnes, Dan B Goldman, Eli Shechtman, and Adam Finkelstein. 2010. Video Tapestries with Continuous Temporal Zoom. *SIGGRAPH* (2010).

Eric P Bennett and Leonard McMillan. 2007. Computational time-lapse video. In *ACM SIGGRAPH 2007 papers*. 102–es.

Daniel Castro, Steven Hickson, Patsorn Sangkloy, Bhavishya Mittal, Sean Dai, James Hays, and Irfan Essa. 2018. Let’s Dance: Learning From Online Dance Videos. In *eprint arXiv:2139179*.

Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. 2019. Everybody dance now. In *Proceedings of the IEEE International Conference on Computer Vision*. 5933–5942.

Yung-Yu Chuang, Aseem Agarwala, Brian Curless, David Salesin, and Richard Szeliski. 2002. Video matting of complex scenes. In *SIGGRAPH*.

Abe Davis and Maneesh Agrawala. 2018. Visual Rhythm and Beat. *ACM Trans. Graph.* 37, 4 (2018), 122–1.

Hao-Shu Fang, Shuqin Xie, Yu-Wing Tai, and Cewu Lu. 2017. RMPE: Regional Multi-person Pose Estimation. In *ICCV*.

Matthieu Fradet, Patrick Pérez, and Philippe Robert. 2008. Semi-automatic Motion Segmentation with Motion Layer Mosaics. In *ECCV*.

Oran Gafni, Lior Wolf, and Yaniv Taigman. 2020. Vid2Game: Controllable Characters Extracted from Real-World Videos. In *ICLR*.

Yossi Gandelsman, Assaf Shocher, and Michal Irani. 2019. “Double-DIP”: Unsupervised Image Decomposition via Coupled Deep-Image-Priors. In *CVPR*.

Dan B. Goldman, Chris Gotterman, Brian Curless, David Salesin, and Steven M. Seitz. 2008. Video Object Annotation, Navigation, and Composition. In *Proceedings of the 21st Annual ACM Symposium on User Interface Software and Technology* (Monterey, CA, USA) (*UIST ’08*). Association for Computing Machinery, New York, NY, USA, 3–12. <https://doi.org/10.1145/1449715.1449719>

Matthias Grundmann, Vivek Kwatra, and Irfan Essa. 2011. Auto-directed video stabilization with robust l1 optimal camera paths. In *CVPR 2011*. IEEE, 225–232.

Rıza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. 2018. Densepose: Dense human pose estimation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 7297–7306.

Qiqi Hou and Feng Liu. 2019. Context-Aware Image Matting for Simultaneous Foreground and Alpha Estimation. In *ICCV*.

Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. 2017. Image-to-Image Translation with Conditional Adversarial Networks. In *CVPR*.

Njegica Jojic and B.J. Frey. 2001. Learning flexible sprites in video layers. In *CVPR*.

Neel Joshi, Wolf Kienzle, Mike Toelle, Matt Uyttendaele, and Michael F Cohen. 2015. Real-time hyperlapse creation via optimal frame selection. *ACM Transactions on Graphics (TOG)* 34, 4 (2015), 1–9.

H. Kim, P. Garrido, A. Tewari, W. Xu, J. Thies, M. Nießner, P. Pérez, C. Richardt, M. Zollhöfer, and C. Theobalt. 2018. Deep Video Portraits. *ACM Transactions on Graphics 2018 (TOG)* (2018).

Diederik Kingma and Jimmy Ba. 2014. Adam: A Method for Stochastic Optimization. *ICLR* (2014).

M. Pawan Kumar, Philip H. S. Torr, and Andrew Zisserman. 2008. Learning Layered Motion Segmentations of Video. *IJCV* (2008).

Shuyue Lan, Rameswar Panda, Qi Zhu, and Amit K Roy-Chowdhury. 2018. FFNet: Video fast-forwarding via reinforcement learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6771–6780.

Jessica Lee, Deva Ramanan, and Rohit Girdhar. 2020. MetaPix: Few-Shot Video Retargeting. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=Sjx1URNKwH>

Yin Li, Jian Sun, and Heung-Yeung Shum. 2005. Video Object Cut and Paste. In *SIGGRAPH*.

Lingjie Liu, Weipeng Xu, Michael Zollhöfer, Hyeonwoo Kim, Florian Bernard, Marc Habermann, Wenping Wang, and Christian Theobalt. 2019. Neural Rendering and Reenactment of Human Actor Videos. *ACM Trans. Graph.* 38, 5, Article 139 (Oct. 2019), 14 pages. <https://doi.org/10.1145/3333002>

Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. 2015. SMPL: A Skinned Multi-Person Linear Model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)* 34, 6 (Oct. 2015), 248:1–248:16.

Ricardo Martin-Brualla, Rohit Pandey, Shuoran Yang, Pavel Pidlipskyi, Jonathan Taylor, Julien Valentin, Sameh Khamis, Philip Davidson, Anastasia Tkach, Peter Lincoln, et al. 2018. LookinGood: Enhancing Performance Capture with Real-Time Neural Re-Rendering. *ACM Trans. Graph.* 37, 6, Article 255 (Dec. 2018), 14 pages. <https://doi.org/10.1145/3272127.3275099>

James McCann, Nancy S Pollard, and Siddhartha Srinivasa. 2006. Physics-based motion retiming. In *Proceedings of the 2006 ACM SIGGRAPH/Eurographics symposium on Computer animation*. Eurographics Association, 205–214.

- Moustafa Meshry, Dan B Goldman, Sameh Khamis, Hugues Hoppe, Rohit Pandey, Noah Snavely, and Ricardo Martin-Brualla. 2019. Neural rendering in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6878–6887.
- Ajay Nandoriya, Elgharib Mohamed, Changil Kim, Mohamed Hefeeda, and Wojciech Matusik. 2017. Video Reflection Removal Through Spatio-Temporal Optimization. In *ICCV*.
- Alasdair Newson, Andrés Almansa, Matthieu Fradet, Yann Gousseau, and Patrick Pérez. 2014. Video Inpainting of Complex Scenes. *SIAM Journal on Imaging Sciences, Society for Industrial and Applied Mathematics* (2014).
- Yair Poleg, Tavi Halperin, Chetan Arora, and Shmuel Peleg. 2015. Egosampling: Fast-forward and stereo for egocentric videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4768–4776.
- Thomas Porter and Tom Duff. 1984. Compositing Digital Images. *SIGGRAPH Comput. Graph.* 18, 3 (Jan. 1984), 253–259. <https://doi.org/10.1145/964965.808606>
- Yael Pritch, Alex Rav-Acha, and Shmuel Peleg. 2008. Nonchronological video synopsis and indexing. *IEEE transactions on pattern analysis and machine intelligence* 30, 11 (2008), 1971–1984.
- Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski. 2011. ORB: An efficient alternative to SIFT or SURF. In *2011 International conference on computer vision*. Ieee, 2564–2571.
- Michel Silva, Washington Ramos, Joao Ferreira, Felipe Chamone, Mario Campos, and Erickson R Nascimento. 2018. A weighted sparse sampling and smoothing frame transition approach for semantic fast-forward first-person videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2383–2392.
- V. Sitzmann, J. Thies, F. Heide, M. Nießner, G. Wetzstein, and Zollhöfer. 2019. Deep-Voxels: Learning Persistent 3D Feature Embeddings. In *Proc. Computer Vision and Pattern Recognition (CVPR)*. IEEE.
- Pratul P. Srinivasan, Richard Tucker, Jonathan T. Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely. 2019. Pushing the Boundaries of View Extrapolation with Multiplane Images. In *CVPR*.
- Justus Thies, Michael Zollhöfer, and Matthias Nießner. 2019. Deferred neural rendering: image synthesis using neural textures. *ACM Trans. Graph.* 38, 4 (2019), 66:1–66:12.
- Giorgio Tomasi, Frans Van Den Berg, and Claus Andersson. 2004. Correlation optimized warping and dynamic time warping as preprocessing methods for chromatographic data. *Journal of Chemometrics: A Journal of the Chemometrics Society* 18, 5 (2004), 231–241.
- Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. 2018. Deep image prior. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 9446–9454.
- John Wang and Edward Adelson. 1994. Representing Moving Images with Layers. *IEEE Trans. on Image Processing* (1994).
- Jue Wang, Pravin Bhat, R. Alex Colburn, Maneesh Agrawala, and Michael F. Cohen. 2005. Interactive Video Cutout. *TOG* (2005).
- Yonatan Wexler, Eli Shechtman, and Michal Irani. 2007. Space-Time Completion of Video. *PAMI* (2007).
- Yuliang Xiu, Jiefeng Li, Haoyu Wang, Yinghong Fang, and Cewu Lu. 2018. Pose Flow: Efficient Online Pose Tracking. In *BMVC*.
- Ning Xu, Brian Price, Scott Cohen, and Thomas Huang. 2017. Deep Image Matting. In *CVPR*.
- Tianfan Xue, Michael Rubinstein, Ce Liu, and William T. Freeman. 2015. A Computational Approach for Obstruction-Free Photography. *ACM Transactions on Graphics (Proc. SIGGRAPH)* 34, 4 (2015).
- Innfarn Yoo, Michel Abdul Massih, Illia Ziamtsov, Raymond Hassan, and Bedrich Benes. 2015. Motion retiming by using bilateral time control surfaces. *Computers & Graphics* 47 (2015), 59–67.
- Feng Zhou, Sing Bing Kang, and Michael F Cohen. 2014. Time-mapping using space-time saliency. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3358–3365.
- Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. 2018. Stereo Magnification: Learning View Synthesis using Multiplane Images. In *SIGGRAPH*.
- Yipin Zhou, Zhaowen Wang, Chen Fang, Trung Bui, and Tamara Berg. 2019. Dance dance generation: Motion transfer for internet videos. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 0–0.
- C. Lawrence Zitnick, Sing Bing Kang, Matthew Uyttendaele, Simon Winder, and Richard Szeliski. 2004. High-quality video view interpolation using a layered representation. *TOG*.

A APPENDIX

A.1 Implementation Details

Neural texture. We use a neural texture with 16 channels, where each body part and the background are represented by an atlas

with 16×16 pixel dimensions. As in the DensePose work, we adopt the SMPL representation, which uses 24 body parts. Thus the dimensions of our neural texture, for a video with N people, is $16 \times 16 \times 16 * (24N - 1)$. Empirically we found that increasing the texture resolution did not improve results.

Training on low-resolution. For each video, we first train the neural renderer and neural texture for 2500 epochs at a video size of 448×256 (or 352×256 for the *Ballroom* sequence only). The input to the neural rendering network is the sampled neural texture concatenated with a map representing the person ID at each pixel (or 0 for background). We apply brightness and spatial jittering to the video frames as data augmentation. Brightness jittering is turned off after 400 epochs to allow the network to faithfully reconstruct the original video. We also employ curriculum learning, where we train only on the easier half of the frames for the first 1k epochs. We rank difficulty by computing the IoU of every pair of layers’ positive trimap regions, using distance of bounding box centers as a tie-breaker. The goal of this curriculum is to withhold the more difficult frames that contain overlapping people, because of the greater ambiguity in learning correlations.

Upsampling low-resolution results. After training the neural rendering module, we obtain a high-resolution result by freezing the trained parameters and training an additional upsampling network. This lightweight upsampling network is trained for 500 epochs with only L1 reconstruction loss, sampling random 256×256 crops for upsampling due to memory constraints. The final output of the upsampling network has dimensions double the size of the low-resolution output. We apply the detail transfer step described in Section 4.4 to these final outputs.

A.2 Network Architectures

In all networks, zero-padding is used to preserve the input shape. ‘bn’ refers to batch normalization. ‘in’ refers to instance normalization. ‘convt’ refers to convolutional transpose. ‘leaky’ refers to leaky relu with slope -0.2. ‘skipk’ refers to a skip connection with layer k . ‘resblock’ denotes a residual block consisting of conv, instance norm, relu, conv, instance norm.

A.2.1 Neural renderer. The neural renderer architecture is a modified pix2pix network [Isola et al. 2017]:

	layer type(s)	out channels	stride	activation
1	4×4 conv	64	2	leaky
2	4×4 conv, bn	128	2	leaky
3	4×4 conv, bn	256	2	leaky
4	4×4 conv, bn	256	2	leaky
5	4×4 conv, bn	256	2	leaky
6	4×4 conv, bn	256	1	leaky
7	4×4 conv, bn	256	1	leaky
8	skip5, 4×4 convt, bn	256	2	relu
9	skip4, 4×4 convt, bn	256	2	relu
10	skip3, 4×4 convt, bn	128	2	relu
11	skip2, 4×4 convt, bn	64	2	relu
12	skip1, 4×4 convt, bn	64	2	relu
13	4×4 conv	4	1	tanh

A.2.2 Upsampling network. The upsampling network predicts a residual image for each bilinearly upsampled RGBA layer predicted by the neural renderer. The network input is the bilinearly upsampled (to the desired output size) concatenation of (1) the predicted low-resolution RGBA, (2) the sampled texture input, and (3) the final feature maps output by the neural renderer preceding the RGBA output layer. The RGBA outputs of the upsampling network are then composited according to σ_t . The upsampling network architecture is as follows:

	layer type(s)	output channels	stride	activation
1	3×3 conv, in	64	1	relu
2	3×3 resblock	64	1	relu
3	3×3 resblock	64	1	relu
4	3×3 resblock	64	1	relu
5	3×3 conv	4	1	none

A.2.3 Keypoints-to-UVs. The keypoint-to-UV network is a fully convolutional network that takes in an RGB image of a skeleton and outputs a UV map of the same size. The architecture is the same as the neural renderer architecture, with the exception of the final layer, which is replaced by two heads: 1) a final convolutional layer with 25 output channels to predict body part and background classification, and 2) a convolutional layer with 48 output channels to regress UV coordinates for each of the 24 body parts. As in the DensePose work [Güler et al. 2018], we train the body part classifier with cross-entropy loss and train the predicted UV coordinates with L1 loss. The regression loss on the UV coordinates is only taken into account for a body part if the pixel lies within the specific part, as defined by the ground-truth UV map.

4

Omnimatte: Associating Objects and Their Effects in Video

This work was presented as an oral presentation at the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2021.

Omnimatte: Associating Objects and Their Effects in Video

Erika Lu^{1,2}

Forrester Cole¹

Tali Dekel^{1,3}

Andrew Zisserman²

William T. Freeman¹

Michael Rubinstein¹

¹Google Research

²University of Oxford

³Weizmann Institute of Science

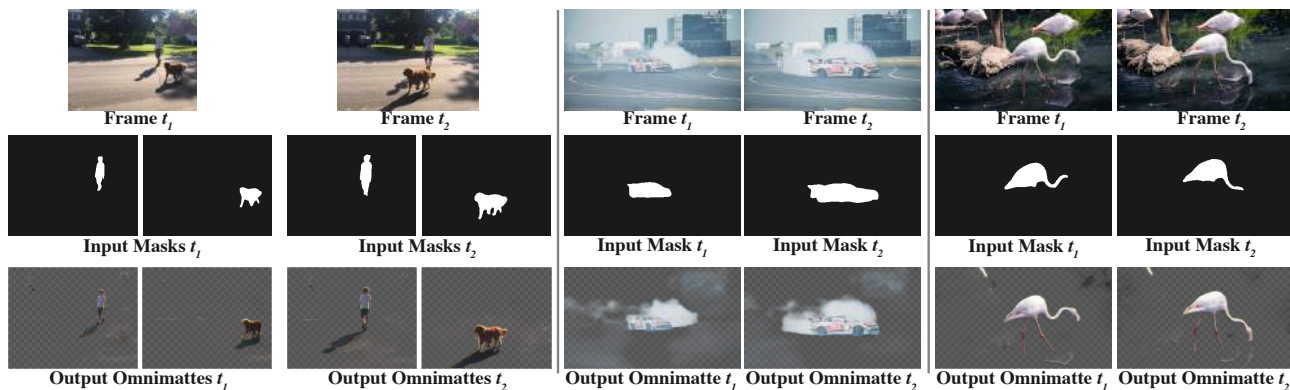


Figure 1. We pose a novel problem: automatically associating subjects in videos with ‘effects’ related to them in the scene. Given an input video (top) and rough masks of subjects of interest (middle), our method estimates an *omnimatte* – an alpha matte and foreground color that includes the subject itself along with all scene elements associated with it (bottom). The associated elements can be other objects attached to the subject or moving with it, or complex effects such as shadows, reflections, smoke, or ripples the subject creates in water.

Abstract

Computer vision is increasingly effective at segmenting objects in images and videos; however, scene effects related to the objects—shadows, reflections, generated smoke, etc.—are typically overlooked. Identifying such scene effects and associating them with the objects producing them is important for improving our fundamental understanding of visual scenes, and can also assist a variety of applications such as removing, duplicating, or enhancing objects in video. In this work, we take a step towards solving this novel problem of automatically associating objects with their effects in video. Given an ordinary video and a rough segmentation mask over time of one or more subjects of interest, we estimate an omnimatte for each subject—an alpha matte and color image that includes the subject along with all its related time-varying scene elements. Our model is trained only on the input video in a self-supervised manner, without any manual labels, and is generic—it produces omnimattes automatically for arbitrary objects and a variety of effects. We show results on real-world videos containing interactions between different types of subjects (cars, animals, people) and complex effects, ranging from semi-transparent elements such as smoke and reflections, to fully opaque effects such as objects attached to the subject.¹

¹Project page: <https://omnimatte.github.io/>

1. Introduction

“And first he will see the shadows best, next the reflections of men and other objects in the water, and then the objects themselves, then he will gaze upon the light of the moon and the stars and the spangled heaven ... Last of all he will be able to see the sun.” – Plato

Is it possible to automatically determine all the effects caused by a subject in a video? Reflect for a moment on the difficulty of the task: a subject, such as a human wandering through a scene, can cast shadows on the floor and distant walls, and be reflected in windows and other surfaces. These ‘effects’ are non-local. However, they are *correlated* with the subject’s shape, motion and, in the case of reflections, appearance.

Tackling this problem is the objective of this paper. More specifically, given an input video and (possibly rough) segmentations over time of subjects of interest in the video, we seek to produce an output opacity matte (alpha matte) for each subject that includes the subject and their effects in the scene (Figure 1). We call this the “*omnimatte*” of the subject. We additionally produce a color background image containing the static background elements in the video. We achieve this by proposing a network and training framework that is able to automatically determine and segment regions that are *correlated* with the given subject (Figure 2). The model is trained in a self-supervised way only on the in-

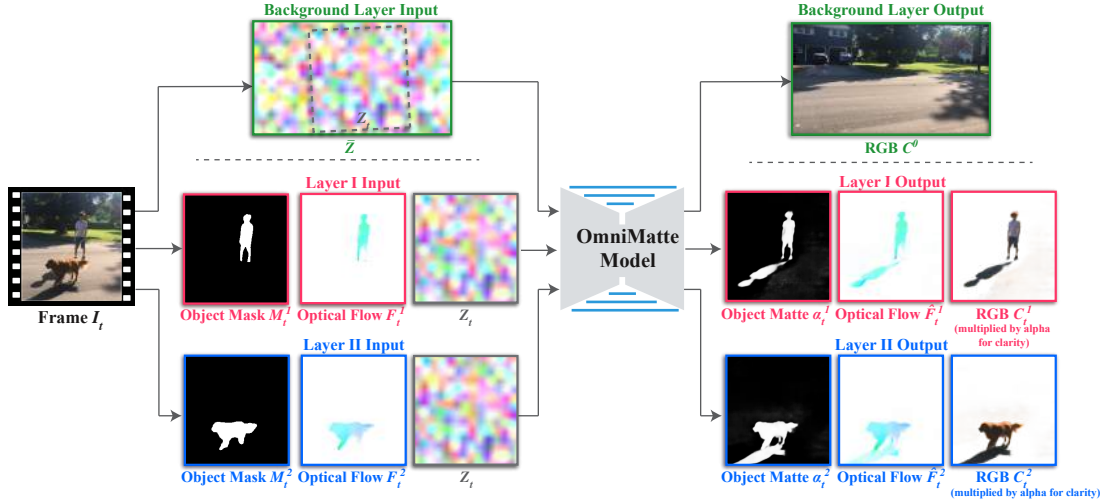


Figure 2. **Estimating omnimattes from video.** The input to the model is an ordinary video with multiple moving objects, and a rough segmentation mask M for each object (left). In a pre-processing step, we compute an optical flow field F between consecutive frames using [30]. For each object, we pass the mask, estimated flow in the object’s region, and a sampled noise image Z_t (representing the background) to our model, producing an omnimatte (color + opacity) and an optical flow field for the object (right). In addition, the model predicts a single background color image for the entire video (top), given a spatial texture noise image \bar{Z} as input. See Sec. 3 for details.

put video, without observing any additional examples. Our solution is inspired by the recent work of Lu et al. [18] that presented a method to decompose a video into a set of human-specific RGBA layers. We generalize this technique to support arbitrary objects, by relying only on binary input masks (no object-specific representation or processing) and incorporating general optical flow to account for motion and frame-to-frame correspondence.

Associating objects with their effects not only improves our fundamental understanding of visual scenes and events captured in video, it can also support a range of applications. Consider for example the problem of removing a person or other types of objects from a video. As is well known, a common error in person removal, e.g., by inpainting, is that a shadow or reflection of the person remains, resulting in a video left with just a ‘shadow of the former self’. The erroneous missing of a reflection is a central plot point in the film ‘Rising Sun’ (1993), and the converse, a lack of reflection, a common trope of vampire movies. The important point is that manipulating an object in a video requires dealing not only with the object; its *effects* in the scene need to be adjusted together with the object in order to create realistic and faithful renditions.

We demonstrate results of inferring omnimattes for different objects such as animals, cars, and people, capturing a variety of complex scene effects including shadows, reflections, dust and smoke. We evaluate the resulting omnimattes qualitatively and quantitatively, and also demonstrate how omnimattes can be useful for video editing applications such as object removal, background replacement, ‘color pop’, and stroboscopic photography.

2. Related Work

Video layer decomposition Our work is inspired by seminal works on layered video decomposition such as [32] and [6]. Layered representations of images and videos

have been applied widely in computer vision and graphics, for example, for inferring occlusion relationships (e.g. [6]), depth (e.g. [41, 29]), mosaics (e.g. [10]), and synthesizing novel views (e.g. [28, 31]). In particular, our work builds on the recent work of Lu et al. [18] that presented a method for decomposing a video into a set of human-specific RGBA layers, where each layer represents a person and their associated scene elements. They take a neural rendering approach and represent the geometry and texture of people explicitly, using a dedicated, human-specific pipeline. We demonstrate that the neural rendering component is in fact unnecessary, and that comparable results can be achieved by providing only rough segmentation masks and flow as input to the network. The result is a simpler and more efficient setup that allows the model to handle arbitrary moving objects. We compare the results with [18] in Sec. 4.5.

Image and video matting Image and video matting traditionally deals with the problem of estimating a foreground layer (color + opacity) and a background color image from a given image or a video (e.g., [5, 33, 16, 7, 39, 14, 27]). The novel problem we propose—estimating an omnimatte—also aims at estimating color + opacity layers from an input video. However, the key fundamental difference is that omnimattes capture not only an object but also all the various scene effects that are correlated with the object. None of the existing matting methods is suitable for performing this task: they cannot handle well entirely semi transparent objects, typically require accurate trimaps that are generated manually, and they are often restricted to estimating two layers (background/foreground). In practice, film production uses manual or semi-automatic rotoscoping to create mattes with such effects [15]. Our method works automatically and generically on natural, ordinary videos that contain arbitrary moving objects and scene effects, and requires only rough object masks (e.g., see the flamingo example in Fig. 1).

Background subtraction Change detection using background subtraction [22, 9] typically does not produce alpha mattes, but binary masks containing all objects and effects (such as shadows). Qian and Sezan’s “difference matting” work [25] attempts to use background subtraction and thresholding to produce a foreground matte with a known background image, but the results are very sensitive to the thresholding value. [27] recently modernized that approach, producing nice quality, continuous-valued alpha mattes but still require a known, clean image of the background. More importantly, background subtraction and difference matting cannot solve the omnimatte problem when the video has *multiple objects with effects*. In such cases it is not enough to detect the effects, each effect must also be associated its subject. We evaluate our method numerically and compare it with background subtraction using a change detection dataset [36] with pixel-level labels for objects and shadows.

Shadow and reflection detection Specialized methods also exist for detecting, removing, or modifying specific types of effects, such as shadows and reflections. For example, [8] acquire a shadow displacement map by waving a stick over different parts of a scene, then use it to synthesize realistic shadows that match an object’s shape. [2, 3] decompose natural videos to remove reflections, shadows, and smoke. More recently, Wang et al. [35] proposed to analyze the motion of people in video and use it to predict depth, occlusion, and lighting/shadow information, to increase realism of 2D object insertion. Our goal in this work is to provide a general technique for inferring all of a subject’s associated effects. However, as shadows are a particularly common effect, we compare our results with a state-of-the-art shadow detector [34] in Sec. 4.1.

Video Effects Although omnimattes are not explicitly optimized for editing, they can facilitate various video editing effects that rely on input object masks, including object removal and video completion (e.g. [11, 37, 19]), object cut-and-paste (e.g. [16, 33]), color pop, and creation of stroboscopic images from video (e.g. [1]). All of these effects can be achieved via simple manipulations of the estimated omnimattes in a post-processing step, or alternatively, by using the omnimattes as input to existing methods such as video completion (e.g. [11]), to save the manual work required for marking the object’s effects. We demonstrate these results in Sec. 4.1.

3. Estimating Omnimattes from Video

The input to our method is an ordinary video of moving objects, and one or more layers of rough segmentation masks that mark the subjects of interest. The output is an *omnimatte* for each input mask layer, consisting of an alpha matte (opacity map) and a color image. The model is trained per-video to reconstruct the input in a self-supervised manner, without observing any additional examples.

To accurately reconstruct the input video, the model must infer all the time-varying effects (e.g. shadows, reflections) from the input object masks, which do *not* represent

those effects. Our goal is to steer the model to place the associated effects in the layer of the subject causing them. Lu, et al. [18] showed that this association can be achieved by showing the network one mask at a time, leveraging the fact that an effect is easier to predict from the object mask most correlated with it. For example, the mask of the person in Figure 1 provides more information about its shadow (more similar to it in shape, in motion) compared to the mask of the dog. Therefore (as shown in [18]) the network tends to learn to predict the person’s shadow from the person’s mask (thus associating it with the correct layer). We build on this training strategy, but design network inputs and losses to encourage this solution for general objects.

3.1. Overview

Figure 2 illustrates our pipeline. Our model is a 2D U-Net [26] that processes the video frame by frame. For each frame, we compute rough object masks using off-the-shelf techniques to mark the major moving objects in the scene. We group the objects into N mask layers $\{M_t^i\}_{i=1}^N$ and define a (possibly time-varying) ordering o_t for the layers. For example, in a scene with a rider, a bicycle, and several people in a crowd, we might group the rider and bicycle into one layer, while grouping the crowd into a second layer. To equip our model with explicit information about *object motion* and frame-to-frame correspondence, we also compute a dense optical flow field, F_t , between each frame and the consecutive frame in the video. This flow field is masked by the input masks M_t^i to provide the network only flow information related to the layer’s subject. We additionally align all frames onto a common coordinate system using homographies, and represent the background as a single unwrapped image on a separate layer.

From this rough yet explicit representation of moving objects, the model has to infer: (i) *omnimattes* – pairs of continuous-valued opacity maps (mattes) and RGB images that capture not only the i^{th} moving object but also all the scene elements that are correlated with it in space and time (e.g., reflections, shadows, attached objects, etc.), (ii) a refined optical flow field for each layer, and (iii) a background RGB image. Formally,

$$\text{Omnimatte}(I_t, H_t, M_t^i, F_t^i) = \mathcal{L}_t = \{\alpha_t^i, C_t^i, \hat{F}_t^i\}, \quad (1)$$

where I_t, H_t, M_t^i, F_t^i are the input video RGB frame, estimated camera homography, the initial input mask, and the pre-computed flow field of the i^{th} object in time t , respectively. α_t^i and C_t^i are the alpha and color buffers of the output omnimatte, and \hat{F}_t^i is the predicted object flow.

The training loss consists of terms on the RGBA outputs (Sec. 3.2) and the predicted flow (Sec. 3.3). The main loss is a reconstruction loss $\mathbf{E}_{\text{rgb-recon}}$, but as reconstruction is underconstrained with multiple layers, we add a sparsity regularization \mathbf{E}_{reg} to the alpha layers and an initialization loss \mathbf{E}_{mask} to the masks, similar to [18]. We encourage the *motion* of the result to match the input by adding a flow-reconstruction loss $\mathbf{E}_{\text{flow-recon}}$ and a temporal consistency term to the alpha mattes $\mathbf{E}_{\text{alpha-warp}}$.

The total loss is:

$$\mathbf{E}_{\text{rgb-recon}} + \lambda_r \mathbf{E}_{\text{reg}} + \lambda_m \mathbf{E}_{\text{mask}} + \mathbf{E}_{\text{flow-recon}} + \lambda_w \mathbf{E}_{\text{alpha-warp}}, \quad (2)$$

where λ_r , λ_m , and λ_w are weighting coefficients (see supplementary material (SM)). As the background is assumed to be static, we factor out camera motion and treat the background with a special, fixed layer (Sec. 3.4).

3.2. RGBA Losses

The main loss in our optimization is a *reconstruction loss*. Formally, we composite the set of estimated layers for each frame and the predicted background layer using standard back-to-front compositing [24], and encourage the composite image to match the original frame:

$$\mathbf{E}_{\text{rgb-recon}} = \frac{1}{T} \sum_t \|I_t - \text{Comp}(\mathcal{L}_t, o_t)\|_1, \quad (3)$$

where $\mathcal{L}_t = \{\alpha_t^i, C_t^i\}_{i=1}^N$ are the predicted layers for frame t , and o_t is the compositing order.

To prevent a trivial solution where a single layer reconstructs the entire frame, we further apply a regularization loss to the α_t^i to encourage them to be spatially sparse. We use a mix of L_1 and an approximate- L_0 :

$$\mathbf{E}_{\text{reg}} = \frac{1}{T} \frac{1}{N} \sum_t \sum_i \gamma \|\alpha_t^i\|_1 + \Phi_0(\alpha_t^i), \quad (4)$$

where $\Phi_0(x) = 2 \cdot \text{Sigmoid}(5x) - 1$ smoothly penalizes non-zero values of the alpha map, and γ controls the relative weight between the terms.

To guide the optimization to convergence from a random initialization, we therefore adopt a ‘‘bootstrap’’ loss to coerce the alpha maps α_t^i to match the input masks M_t^i :

$$\mathbf{E}_{\text{mask}} = \frac{1}{T} \frac{1}{N} \sum_t \sum_i \|d_t^i \odot (M_t^i - \alpha_t^i)\|_2 \quad (5)$$

where $d_t^i = 1 - \text{dilate}(M_t^i) + M_t^i$ is a boundary erosion mask to turn off the loss near the mask boundary, and \odot is element-wise product. This loss is turned off after its value reaches a fixed threshold (see SM).

3.3. Flow Losses

Our model additionally predicts a set of *flow layers*. Predicting flow layers serves as an auxiliary task that injects information about motion to our model and improves our decomposition (as demonstrated by our experiments). To achieve that we apply a *flow reconstruction loss* and a *photometric warping loss* defined below:

$$\mathbf{E}_{\text{flow-recon}} = \frac{1}{T} \sum_t W_t \cdot \|F_t - \text{Comp}(\mathcal{F}_t, o_t)\|_1, \quad (6)$$

where $\mathcal{F}_t = \{\hat{F}_t^i\}$ is the set of predicted flow layers, F_t is the original, pre-computed flow, and W_t is a spatial weighting map that lowers the impact of pixels with inaccurate flow. W_t is computed based on standard left-right flow consistency error and photometric warping error (see full details in SM).

We additionally encourage temporal consistency within layers using an *alpha warping loss*:

$$\mathbf{E}_{\text{alpha-warp}} = \frac{1}{T} \frac{1}{N} \sum_t \sum_i \|\alpha_t^i - \alpha_{wt}^i\|_1, \quad (7)$$

where $\alpha_{wt}^i = \text{Warp}(\alpha_{t+1}^i, \mathcal{F}_t^i)$ is the alpha for layer i at time $t + 1$ warped to time t using the predicted flow.

3.4. Camera Motion and Background

We assume the background scene is stationary and camera motion can be modeled by a time-varying homography from an unwrapped ‘‘canvas’’ image, as in [32]. The homographies H_t from frame t to the canvas are estimated via feature tracking (using [12]) on the original RGB video frames and are held fixed. For input to the network, the background canvas is represented by a single spatial noise image \bar{Z} (see Fig. 2). The background color layers C_t^0 are produced by feeding \bar{Z} through the network to form a static color image C^0 , which is then sampled using H_t^{-1} to form time-varying background images $\{C_t^0\}$.

To make the foreground layers aware of the camera motion, the input mask layers M_t^i are concatenated with a noise image that tracks the camera. The background noise image \bar{Z} is sampled using H_t^{-1} to form time-varying noise images $\{Z_t\}$. This is a similar approach to Lu, et al. [18], but in our case the noise image is not trainable.

Minor stabilization errors, as well as exposure changes, vignetting, and radial distortion, usually cause slight changes in appearance even for a stationary background. If the background is assumed to be entirely static, these subtle shifts in appearance will show up as noise in our omnimatte. Such effects, however, tend to have low spatial and temporal frequency relative to the subject’s effects and can be safely captured by applying a *refinement warp* to the background layer. The refinement warp consists of a spatially and temporally coarse grid-based warp. We additionally apply a grid-based brightness adjustment to the final composite $\text{Comp}(\mathcal{L}_t, o_t)$. The parameters of the warp and brightness adjustment are optimized together with the network parameters (see SM for additional details).

3.5. Implementation Details

For all our results, we used Mask R-CNN [13] to segment the input objects, and STM [20] (a video object segmenter trained on the DAVIS dataset [23]) to track objects across frames. Optical flow between consecutive frames was computed using RAFT [30]. When dynamic background elements such as tree branches are present, we use panoptic segmentation [38] to segment them and treat the segment as additional objects. To increase the detail of the color buffers C_t^i , we apply a similar detail-transfer technique to Lu, et al. [18]. See SM for training details.

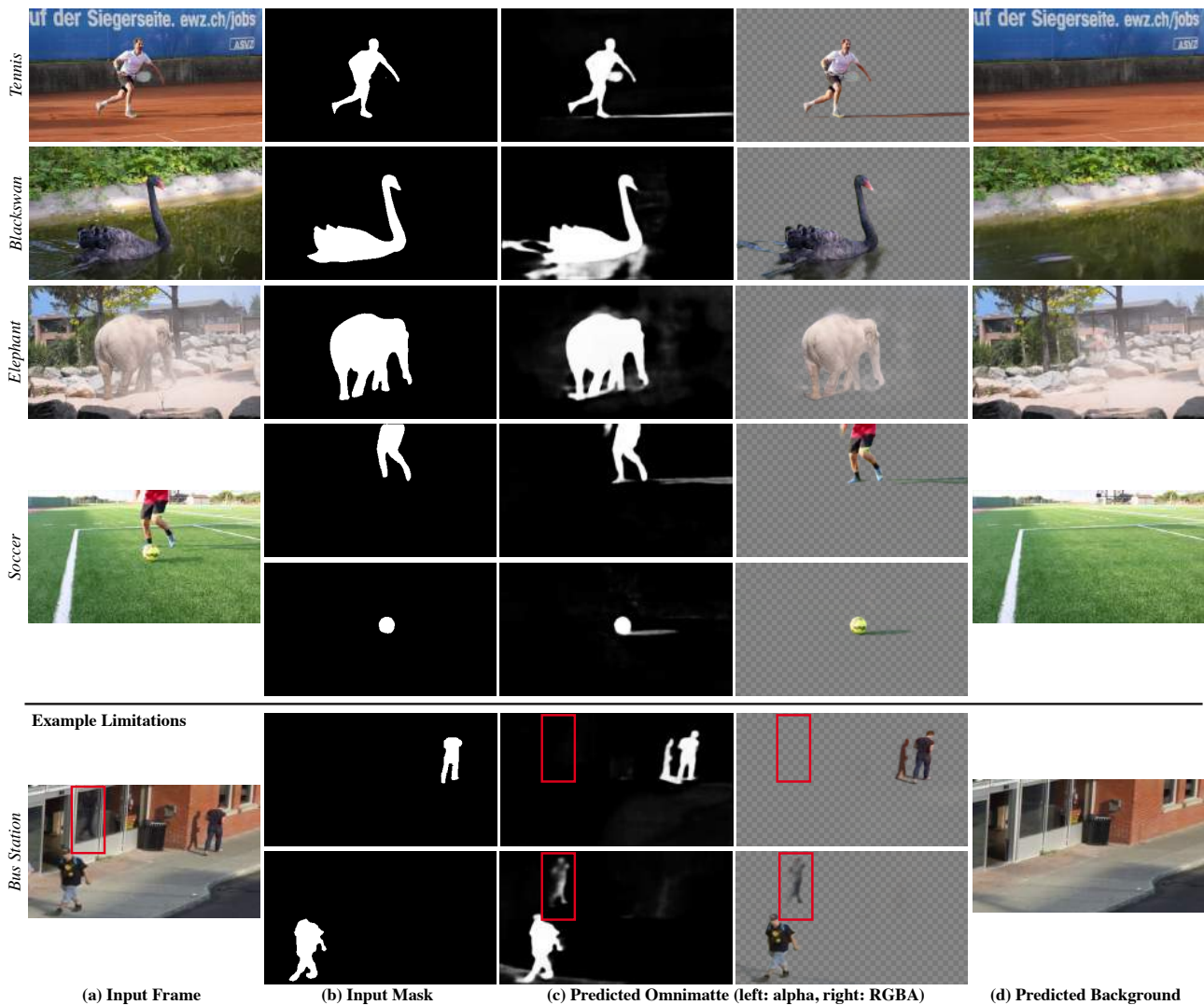


Figure 3. **Results on natural videos.** For each example, we show: (a) input frame; (b) input mask(s) computed by Mask R-CNN [13]; (c) our resulting omnimatte (left: alpha matte, right: RGBA); (d) our estimated background layer. The bottom example (*Bus Station*) shows a failure case: while the shadows are correctly associated with the people, the reflection cast on the window by the person in the top-right corner (marked by the red rectangle) is mistakenly grouped with the person in the bottom-left corner.

4. Results

4.1. Qualitative examples on real videos

Figure 3 shows examples of our estimated omnimattes on a variety of real-world videos from DAVIS [23], CDW-2014 [36] (see Sec. 4.4), and videos downloaded from YouTube. These examples span a wide range of dynamic subjects (e.g., people, animals or general moving objects such as a soccer ball), performing complex actions and generating various scene effects including shadows, reflections, water ripples, dust and smoke. None of the input object masks include these effects (see Fig. 3(b)).

As seen in Fig. 3(c-d) top, our method successfully associates the subjects with the scene effects that are related to them. In *Blackswan*, the omnimatte of the swan captures its reflection and the water ripples it causes. In *Elephant*, our omnimatte captures the semi-transparent cloud of dust

sprayed by the elephant, as well as the shadow the elephant casts on the ground. In *Tennis*, the running player casts thin shadows, which the omnimatte correctly separates from the shadows in the background. Additionally, although the player’s racket is not included in the input mask (b), it is reconstructed in our omnimatte result; this demonstrates our model’s ability to reconstruct objects that are attached to the main subject even when given incomplete input masks.

In *Soccer*, we show a two-subject example where our model estimates a separate omnimatte for each of the subjects: a person (top row) and a soccerball (bottom row). Our model successfully separates the person’s shadow from that of the soccerball up until the final few frames of the video, where part of the person’s shadow appears in the soccerball’s omnimatte (full video in SM).

Another two-subject example is shown in *Bus Station* where two people walk away from each other. Our model

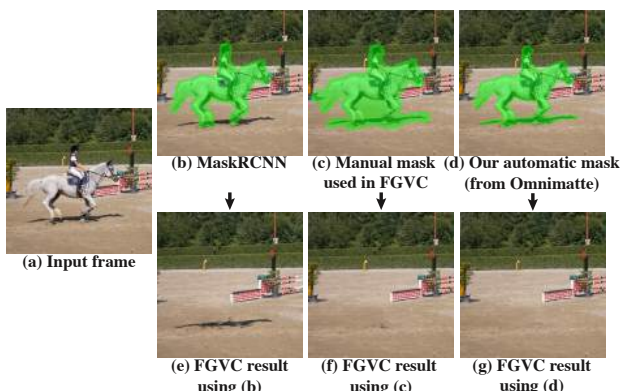


Figure 4. **Omnimattes as input to state-of-the-art object removal.** Results by FGVC [11] using different types of input masks: (b) Raw masks from MaskRCNN do not capture shadows, and produce unrealistic results (e). By using manually annotated masks that include the shadow (c), both the horse and the shadow are removed (f). (d) binary mask *automatically derived from omnimate* produces comparable result (g) when inputted to FGVC.

correctly associates each person with their shadow in this challenging case. However, the reflection cast on the window by the person on the right (top) is incorrectly placed in the left person’s layer (bottom). The challenge of this scene lies in both the spatial proximity of the reflection to the incorrect person, and the similar motions (both people in the scene are moving consistently at the same speed). Lu, et al. [18] showed that layers tend to ‘grab’ spatially proximal effects; in this case, the reflection is actually closer to the person who is *not* casting the reflection. [18] additionally showed that correlated motions are grouped in the same layers; as both people are walking in synchronization, the network places the reflection in the incorrect person’s layer.

4.2. Object Removal

Our method can be applied to remove a dynamic object from a video by either: (i) binarizing our omnimate and using it as input to a separate video-completion method such as FGVC [11], or by (ii) simply excluding the object’s omnimate layer from our reconstruction.

As shown in Fig. 4(b,e), removing an object but not its correlated effects produces an unrealistic result (object removed but its shadow remains). Typically such effects are manually annotated to create a conservative binary mask of the regions to remove (Fig. 4(c)). To show that an omnimate can replace manual editing, we derive a binary mask by thresholding our soft alpha at 0.25 and dilating by 20 pixels, and inputting it to FGVC [11]. Fig. 4(c) shows both the horse and its shadow are removed, demonstrating that our derived mask is comparable to a manually annotated mask.

Fig. 5 shows a comparison between omnimate removal (approach (ii) above) and FGVC using manual masks. In the *flamingo* example, our method removes not only the flamingo but also its reflection in the water beneath it. FGVC relies on a mask that does not include the reflection, thus the reflection remains intact in their result. Our omnimates bypass the need to manually label such semi-

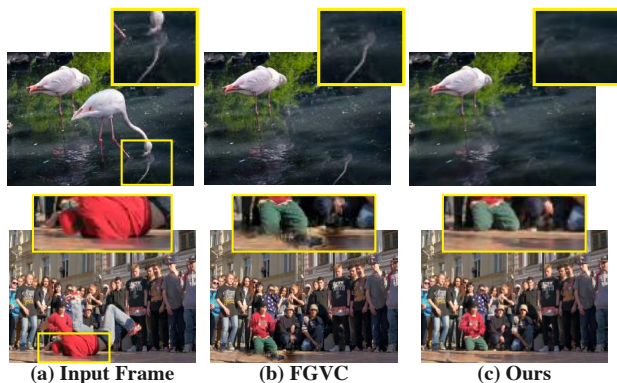


Figure 5. **Direct omnimate-based removal.** For each input frame (a) we show the result of removing a foreground object by excluding its omnimate from the reconstruction (c), compared to FGVC [11] (b).

transparent effects. In the *breakdance* example, both the crowd and the dancer are moving. To handle this case, we assign an omnimate to the dancer and a separate single omnimate to the crowd. Fig. 5(c) shows the crowd omnimate composited with the background layer. The FGVC result (b) shows artifacts on the ground where the dancer is removed, whereas our result is seamless and realistic.

4.3. Comparison with Shadow Detection

We show qualitative comparisons with a recent state-of-the-art shadow detection method, ISD [34], a deep-learning based method that takes an RGB image as input and produces segments for object-shadow pairs. ISD integrates a MaskRCNN-like object detection stage (Detectron2 [38]), hence it does not require or allow an input mask.

Figure 7 compares our result with ISD on two challenging scenes, where a person casts a shadow onto another object (a bench), and where a person’s shadow is occluded by another object (a dog). Our method successfully handles and outperforms ISD in both cases. Occlusions and shadows cast on other objects present particularly difficult cases for purely data-driven methods such as ISD, since the appearance of the shadow depends on the relative configuration of multiple objects in the scene, presenting a combinatorial explosion of scenarios for training. In contrast, our method analyzes and leverages space-time information throughout the entire video to perform these complex object-effects associations.

4.4. Comparison with Background Subtraction

We quantitatively evaluate our approach on the task of background subtraction using a change detection dataset, CDW-2014 [36], which has ground-truth pixel-level labels for objects and hard shadows. We selected a subset of videos that contain objects and their shadows (see Fig. 8 for sample frames and labels). We manually excluded “bad weather” and “low framerate” categories to avoid low quality videos, and selected short clips of up to 5 moving objects. The selected subset contains 12 clips, each with 40 - 115 frames, for approximately 950 frames in total. While the background subtraction task requires only sep-



Figure 6. **Video editing with Omnimattes.** Effects such as “color pop” (left; subject in color, background in grayscale), background replacement (center), and stroboscopic photography (right) all benefit from capturing the subjects’ associated effects with an Omnimatte. Note the color pop present in the flamingo’s reflection and the correct placement of shadows in the background replacement and stroboscopic photograph examples. See the SM for the full details and before/after videos.

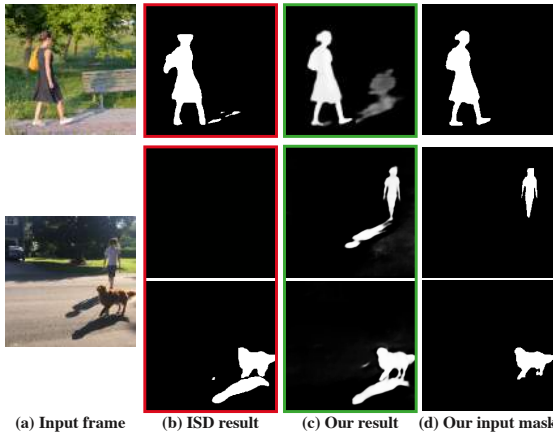


Figure 7. **Comparison with shadow detection.** (b) Results produced by ISD [34], a recent state of the art, single-image shadow detection method, and (c) our results when using (d) MaskRCNN masks as input. See SM for additional comparisons.

arating foreground from background, we demonstrate the additional capabilities of our method by also segmenting the effects for individual object instances.

We convert our soft omnimattes into a single, hard segmentation mask using a fixed threshold value and report the Jaccard index (\mathcal{J}) and Boundary measure (\mathcal{F}) [21] in Table 1. We compare with two top-performing methods on CDW-2014, FgSegNet [17] and BSPVGAN [40], which were trained on subsets of CDW-2014. Our method outperforms FgSegNet and matches the performance of BSPVGAN, despite not being trained supervised on CDW-2014.

4.5. Comparison with Layered Neural Rendering

Fig. 10 shows a qualitative comparison with the human-specific, layered neural rendering method by Lu et al. [18]. In [18], people are parameterized explicitly using per-frame UV maps that represent each individual’s geometry, and a per-person trainable texture map that represents appearance. Instead, we use binary masks and pre-computed optical flow to represent object regions (see Sec. 3). For comparison, we used binary masks extracted from their UV maps.

In both examples, our method achieves comparable results to [18], successfully capturing the trampoline deformations, shadows and reflections, yet with a generic, much simpler input. Note that the input masks derived from the UV maps provided by [18] represent the full body of a person even if they are occluded in the original frame. This

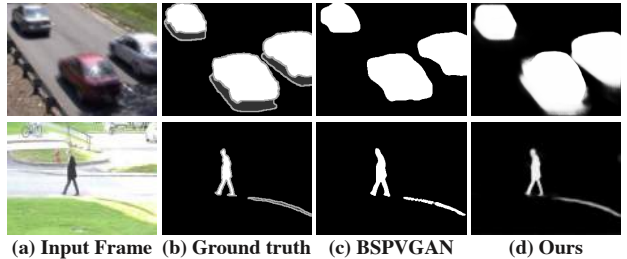


Figure 8. **Comparison with background subtraction.** We used selected videos from the CDW-2014 change detection dataset [36] (examples input frames in (a)), with ground truth, manually segmented objects and shadows (b, white pixels = moving objects, dark gray pixels = shadows, light gray pixels = ‘unknown’, typically at boundaries). (c) Result by a top-performing (on this dataset) background subtraction method [40]. (d) Our result (alpha mattes of estimated omnimattes). Numerical experiments are summarized in Table 1. More results can be found in the SM.

Method	$\mathcal{J} \& \mathcal{F}(\text{Mean}) \uparrow$	$\mathcal{J}(\text{Mean}) \uparrow$	$\mathcal{F}(\text{Mean}) \uparrow$
FgSegNet [17]	0.675	0.631	0.719
BSPVGAN [40]	0.756	0.718	0.793
Ours	0.754	0.711	0.797

Table 1. We compare our method to the two top-performing methods on CDW-2014 [36]. We report the Jaccard index (\mathcal{J}) and Boundary measure (\mathcal{F}) on a subset of the data that includes objects and their shadows. Our method performs at par or better than the two background subtraction methods.

allows our model to inpaint object regions and scene effects that are occluded in some frames but visible in others, as in [18]. The ability of our model to inpaint occluded regions even when using incomplete masks is also evident in the person-dog example in Fig. 1, where the person and their shadow are reconstructed in our omnimatte. However, we note that in cases where the input mask is substantially occluded, the output omnimatte will show occlusion as well; thus in order to deal with large occlusions, a full-object mask should be inputted to the model, as done in [18].

4.6. Additional Video Editing Effects

The additional information present in an omnimatte compared to a standard matte that includes only the subject allow simple creation of various video effects such as color pop, background replacement, or object duplication (Fig. 6). Previously, creating these effects for videos containing shadows or reflections required extensive manual editing effort. Since the omnimatte is a standard RGBA im-

Figure 9. **Ablations.** We ablate several components of our method: **Left: (a-b)** sample input frame and our result using our full method, and **(c)** removing the brightness adjustment, and **(d)** removing the background offset (Section 3.4). **Right: we show two examples comparing our method with and without the flow component (f-g).**

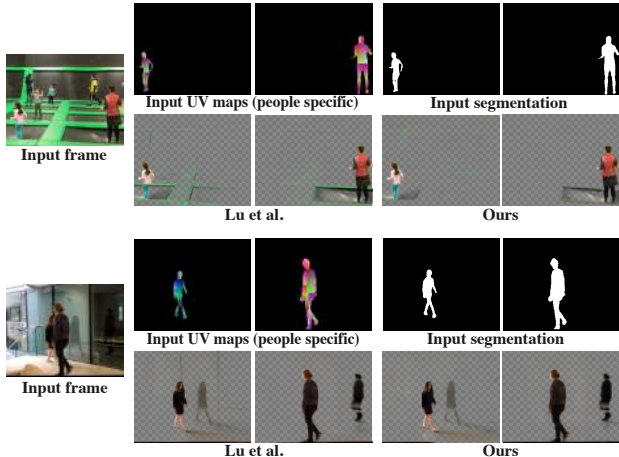
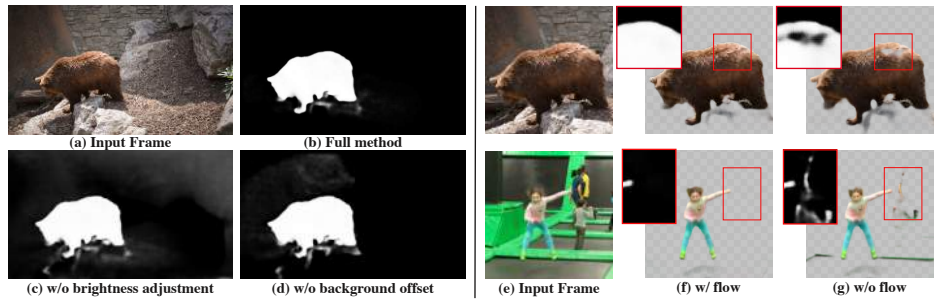


Figure 10. **Comparison with Lu et al. [18].** We achieve comparable results to [18] using just binary masks instead of the people-specific UV maps used in [18]. These binary segmentation masks are easier to obtain and are *general* – allowing our method to support arbitrary objects ([18] is applicable just to people). Notice how our omnimattes capture trampoline deformation well (top two rows), and reflections in the glass (bottom two rows).

age, these edits may be applied using standard video editing software. The omnimattes for color pop and background replacement were used unchanged, the horse jump alpha matte was adjusted with a linear contrast ramp. Please see SM for full details on creating these effects.

4.7. Ablations

In Fig. 9 we ablate several components of our method. Removing the brightness adjustment (c) and removing the background offset (d) both result in undesirable nonzero alpha values in the bear’s omnimatte, due to lighting changes, vignetting, and homography inaccuracies that break the static background assumption (Sec. 3.4). Including both components results in a clean alpha matte (b).

We ablate the flow component of our model by removing both flow inputs and flow losses (Sec. 3.3), and show results in (g). In the top row, part of the bear is missing from the omnimatte, and in the bottom row, the person’s omnimatte incorrectly contains parts of other people. In contrast, our full model (f) has a complete bear and a clean person omnimatte. These examples show that providing the model with motion information (flow) allows it to better associate scene elements with the correct objects, and prevents holes appearing in the foreground object’s alpha matte.

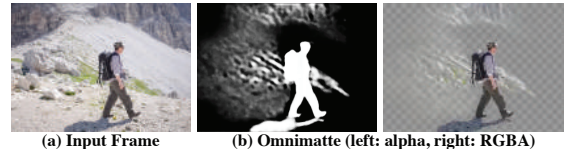


Figure 11. **Failure case due to incorrect camera registration.** When the background motion cannot be accurately represented by a homography (in this case due to a significant depth variation in the scene), the predicted omnimatte may contain regions of the background to compensate for the registration inaccuracies.

4.8. Limitations

While our method allows for small deviations from a static background via smooth, coarse geometric and photometric offsets, when the homographies do not accurately represent the background, the omnimattes must correct for these errors by including background elements (e.g. rocks and grass, Fig. 11). Conversely, we cannot separate objects or effects that remain entirely stationary relative to the background throughout the video. These issues could be addressed by building a background representation that explicitly models the 3D structure of the scene (e.g. [4]).

Finally, we observed that different random initializations of the network’s weights may occasionally lead to different, sometimes undesirable, solutions (see supplemental material for visualization). We speculate that more reliable convergence could be obtained by further optimizing the order in which frames are introduced to the model.

5. Conclusion

We have posed a new problem: from an input video with one or more segmented moving subjects, we produce an *omnimatte* for each subject – an opacity map and color image that includes the subject itself along with the visual effects related to it. These effects can be reflections of the subject, shadows they cast, or attached objects. We have proposed a network and training framework for solving this new problem, and have demonstrated omnimattes produced automatically for real-world videos with a variety of objects and associated effects. We have also shown how omnimattes can support a variety of video editing applications.

Acknowledgements. This work was supported in part by an Oxford-Google DeepMind Graduate Scholarship and a Royal Society Research Professorship. We thank Weidi Xie for assisting with object removal baselines.

References

- [1] Aseem Agarwala, Mira Dontcheva, Maneesh Agrawala, Steven Drucker, Alex Colburn, Brian Curless, David Salesin, and Michael Cohen. Interactive digital photomontage. In *ACM SIGGRAPH 2004 Papers*, pages 294–302. 2004.
- [2] Jean-Baptiste Alayrac, João Carreira, and Andrew Zisserman. The visual centrifuge: Model-free layered video representations. In *CVPR*, 2019.
- [3] Jean-Baptiste Alayrac, Joao Carreira, Relja Arandjelovic, and Andrew Zisserman. Controllable attention for structured layered video decomposition. In *ICCV*, 2019.
- [4] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. In *ECCV*, 2020.
- [5] Xue Bai, Jue Wang, David Simons, and Guillermo Sapiro. Video snapchat: robust video object cutout using localized classifiers. *TOG*, 2009.
- [6] Gabriel J Brostow and Irfan A Essa. Motion based decomposing of video. In *ICCV*, 1999.
- [7] Yung-Yu Chuang, Aseem Agarwala, Brian Curless, David Salesin, and Richard Szeliski. Video matting of complex scenes. In *SIGGRAPH*, 2002.
- [8] Yung-Yu Chuang, Dan B Goldman, Brian Curless, David H Salesin, and Richard Szeliski. Shadow matting and compositing. In *SIGGRAPH*, 2003.
- [9] Ahmed Elgammal, David Harwood, and Larry Davis. Non-parametric model for background subtraction. In *ECCV*, 2000.
- [10] Matthieu Fradet, Patrick Pérez, and Philippe Robert. Semi-automatic motion segmentation with motion layer mosaics. In *ECCV*, 2008.
- [11] Chen Gao, Ayush Saraf, Jia-Bin Huang, and Johannes Kopf. Flow-edge guided video completion. In *ECCV*, 2020.
- [12] Matthias Grundmann, Vivek Kwatra, and Irfan Essa. Auto-directed video stabilization with robust 11 optimal camera paths. In *CVPR*, 2011.
- [13] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask R-CNN. In *ICCV*, 2017.
- [14] Qiqi Hou and Feng Liu. Context-aware image matting for simultaneous foreground and alpha estimation. In *ICCV*, 2019.
- [15] Wenbin Li, Fabio Viola, Jonathan Starck, Gabriel J Brostow, and Neill DF Campbell. Roto++ accelerating professional rotoscoping using shape manifolds. *ACM Transactions on Graphics (TOG)*, 35(4):1–15, 2016.
- [16] Yin Li, Jian Sun, and Heung-Yeung Shum. Video object cut and paste. In *SIGGRAPH*, 2005.
- [17] Long Ang Lim and Hacer Yalim Keles. Learning multi-scale features for foreground segmentation. *arXiv preprint arXiv:1808.01477*, 2018.
- [18] Erika Lu, Forrester Cole, Tali Dekel, Weidi Xie, Andrew Zisserman, David Salesin, William T Freeman, and Michael Rubinstein. Layered neural rendering for retiming people in video. In *SIGGRAPH Asia*, 2020.
- [19] Alasdair Newson, Andrés Almansa, Matthieu Fradet, Yann Gousseau, and Patrick Pérez. Video inpainting of complex scenes. *SIAM Journal on Imaging Sciences*, 2014.
- [20] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *ICCV*, 2019.
- [21] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *CVPR*, 2016.
- [22] Massimo Piccardi. Background subtraction techniques: a review. In *2004 IEEE International Conference on Systems, Man and Cybernetics (IEEE Cat. No. 04CH37583)*, 2004.
- [23] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alexander Sorkine-Hornung, and Luc Van Gool. The 2017 DAVIS challenge on video object segmentation. *arXiv:1704.00675*, 2017.
- [24] Thomas Porter and Tom Duff. Compositing digital images. *SIGGRAPH Comput. Graph.*, 18(3):253–259, Jan. 1984.
- [25] Richard J Qian and M Ibrahim Sezan. Video background replacement without a blue screen. In *Proceedings 1999 International Conference on Image Processing (Cat. 99CH36348)*, volume 4, pages 143–146. IEEE, 1999.
- [26] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. *MICCAI*, 2015.
- [27] Soumyadip Sengupta, Vivek Jayaram, Brian Curless, Steven M. Seitz, and Ira Kemelmacher-Shlizerman. Background matting: The world is your green screen. In *CVPR*, 2020.
- [28] Jonathan Shade, Steven Gortler, Li-wei He, and Richard Szeliski. Layered depth images. In *Proceedings of the 25th annual conference on Computer graphics and interactive techniques*, pages 231–242, 1998.
- [29] Pratul P. Srinivasan, Richard Tucker, Jonathan T. Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely. Pushing the boundaries of view extrapolation with multiplane images. In *CVPR*, 2019.
- [30] Zachary Teed and Jia Deng. RAFT: Recurrent all-pairs field transforms for optical flow. In *ECCV*, 2020.
- [31] Shubham Tulsiani, Richard Tucker, and Noah Snavely. Layer-structured 3D scene inference via view synthesis. In *ECCV*, 2018.
- [32] John Wang and Edward Adelson. Representing moving images with layers. *IEEE Trans. on Image Processing*, 1994.
- [33] Jue Wang, Pravin Bhat, R. Alex Colburn, Maneesh Agrawala, and Michael F. Cohen. Interactive video cutout. *TOG*, 2005.
- [34] Tianyu Wang, Xiaowei Hu, Qiong Wang, Pheng-Ann Heng, and Chi-Wing Fu. Instance shadow detection. In *CVPR*, 2020.
- [35] Yifan Wang, Brian L. Curless, and Steven M. Seitz. People as scene probes. In *ECCV*, 2020.
- [36] Yi Wang, Pierre-Marc Jodoin, Fatih Porikli, Janusz Konrad, Yannick Benezeth, and Prakash Ishwar. CDnet 2014: An expanded change detection benchmark dataset. In *CVPR Workshop*, 2014.
- [37] Yonatan Wexler, Eli Shechtman, and Michal Irani. Space-time completion of video. *PAMI*, 2007.
- [38] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. <https://github.com/facebookresearch/detectron2>, 2019.
- [39] Ning Xu, Brian Price, Scott Cohen, and Thomas Huang. Deep image matting. In *CVPR*, 2017.
- [40] Wenbo Zheng, Kunfeng Wang, and Fei-Yue Wang. A novel background subtraction algorithm based on parallel vision and bayesian GANs. *Neurocomputing*, 2020.

- [41] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *SIGGRAPH*, 2018.

Omnimatte: Associating Objects and Their Effects in Video

Supplementary Material

1. Formulation Details

1.1. Flow Reconstruction Weight (Section 3.3, Eq 6)

The confidence at each pixel p is defined as:

$$W_t(p) = W_t^{lr}(p) \cdot W_t^p(p) \cdot M_t(p), \quad (1)$$

where W_t^{lr} is computed using a standard left-right flow consistency error by: $W_t^{lr}(p) = \max(1 - e_t^{lr}(p), 0)$, and e_t^{lr} is the forward-backward flow error. W_t^p measures photometric error and is given by: $W_t^p = \mathbb{1}_{e_p < \beta}$, where $e_p = \|\text{Warp}(I_{t+1}; F_{t,t+1}) - I_t\|_1$ measures the photometric difference between I_t and I_{t+1} when backward warped using the flow, and $\beta = 20$. Finally, we mask the regions outside of the object mask M_t since semi-transparent effects such as shadows tend to have inaccurate flow.

1.2. Warp and Brightness Adjustment (Section 3.4)

To compensate for minor camera stabilization errors, we apply the learnable warp at each pixel \mathbf{x} of the background layer:

$$\hat{C}_t^0(\mathbf{x}) = C_t^0(\mathbf{G}_w[t, x_v, x_u]) \quad (2)$$

where \mathbf{G}_w is a $n/10 \times 4 \times 7 \times 2$ grid of offset vectors, n is the number of frames in the video, and $[\cdot, \cdot, \cdot]$ denotes trilinear filtering. We additionally apply a learnable coarse brightness scaling to the final composite:

$$\hat{\text{Comp}}(\mathcal{L}_t, o_t) = \text{Comp}(\mathcal{L}_t, o_t) \cdot \mathbf{G}_b[t, x_v, x_u] \quad (3)$$

where \mathbf{G}_b is a $n/10 \times 4 \times 7$ grid of brightness coefficients. The values of \mathbf{G}_w and \mathbf{G}_b are optimized together with the network parameters.

The effect of this adjustment layer can be seen in Fig. 1. Adding the background adjustment significantly increases the sparsity of the alpha mattes for the same reconstruction quality.

1.3. RGBA Detail transfer (Section 3.5)

We adopt the same detail transfer step as in Lu, et al. [6] to produce high-resolution omnimattes by transferring detail from the original frame to the CNN outputs. We first compute the residual between the CNN output and the original frame, and determine the amount of the residual to transfer to each RGBA layer using the transmittance map τ_t^i for

layer i at time t :

$$\tau_t^i = 1.0 - \text{Comp}_\alpha(\mathcal{L}_t \setminus \{L_t^j \mid j < i\}, o_t \setminus \{j \mid j < i\}) \quad (4)$$

where Comp_α denotes the alpha channel of the composite produced by the network. The final layer colors with the transferred detail are:

$$C_t^i = \text{Cnr}_t^i + \tau_t^i(I_t - \text{Comp}(\mathcal{L}_t, o_t)) \quad (5)$$

where Cnr is the color produced by the network.

2. Implementation Details

Network Architecture We adopt the same network architecture as in Lu, et al. [6], with the exception of replacing batchnorm with instance norm:

	layer type(s)	channels	stride	activation
1	conv	64	2	leaky
2	conv, IN	128	2	leaky
3	conv, IN	256	2	leaky
4	conv, IN	256	2	leaky
5	conv, IN	256	2	leaky
6	conv, IN	256	1	leaky
7	conv, IN	256	1	leaky
8	skip5, convt, IN	256	2	relu
9	skip4, convt, IN	256	2	relu
10	skip3, convt, IN	128	2	relu
11	skip2, convt, IN	64	2	relu
12	skip1, convt, IN	64	2	relu
13	conv	4	1	tanh

‘IN’ refers to instance normalization, ‘convt’ refers to convolutional transpose. All convolutions are 4×4 . Additionally, we have a convolutional layer following layer 12 which outputs 2 channels for optical flow.

Training Details We implement our network in JAX [1] and Haiku [3]. We optimize our full objective (Eq. 2, Section 3.1), with relative weights: $\lambda_r = .005$, $\lambda_m = 50$ until E_{mask} falls below 0.05, and is 0 afterward, and $\lambda_w = .005$. We use the Adam optimizer [4] with an initial learning rate of $1e-3$. We train each video with a batch size of 32 on

Google Cloud v3-8 TPU for 2000 epochs. We resize all videos to 256×448 spatial resolution. Training the ‘Bear’ sequence from the DAVIS dataset [7], which has 82 frames and 1 output omnimatte, takes 2 hours.

2.1. Video effects details (Section 4.6)

All the video effects were created in *postprocessing* in NUKE [5], a standard video compositing package. That is, we produce the effects by editing our output omnimattes (and using the camera homographies), whereas the omnimatte network model itself is not used.

Color Pop: to create the “color pop” effect, we create two versions of the original video, one with color desaturated and one with color amplified. We blend the amplified version over the desaturated version using the alpha matte from the foreground layer.

Background Replacement: replacing the background while preserving camera motion is accomplished by treating the new image as a new extended canvas. Thus, we create a new time-varying background frames C_t^0 by applying the original camera homographies.

Stroboscopic Photograph: to create the background of the stroboscopic photograph, we apply the inverse camera homographies to the background color images C_t^0 and accumulate them into the canvas space with over compositing. This step creates a clean background canvas without the foreground objects. We then apply the same inverse transformations to the foreground layers and composite them on top. For the “horsejump-low” photo, we picked 15 frame intervals.

References

- [1] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, and Skye Wanderman-Milne. JAX: composable transformations of Python+NumPy programs, 2018.
- [2] Chen Gao, Ayush Saraf, Jia-Bin Huang, and Johannes Kopf. Flow-edge guided video completion. In *ECCV*, 2020.
- [3] Tom Hennigan, Trevor Cai, Tamara Norman, and Igor Babuschkin. Haiku: Sonnet for JAX, 2020.
- [4] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2014.
- [5] The Foundry Visionmongers Ltd. Nuke, 2018.
- [6] Erika Lu, Forrester Cole, Tali Dekel, Weidi Xie, Andrew Zisserman, David Salesin, William T Freeman, and Michael Rubinstein. Layered neural rendering for retiming people in video. In *SIGGRAPH Asia*, 2020.
- [7] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alexander Sorkine-Hornung, and Luc Van Gool. The 2017 DAVIS challenge on video object segmentation. *arXiv:1704.00675*, 2017.
- [8] Tianyu Wang, Xiaowei Hu, Qiong Wang, Pheng-Ann Heng, and Chi-Wing Fu. Instance shadow detection. In *CVPR*, 2020.
- [9] Yi Wang, Pierre-Marc Jodoin, Fatih Porikli, Janusz Konrad, Yannick Benezeth, and Prakash Ishwar. CDnet 2014: An expanded change detection benchmark dataset. In *CVPR Workshop*, 2014.
- [10] Wenbo Zheng, Kunfeng Wang, and Fei-Yue Wang. A novel background subtraction algorithm based on parallel vision and bayesian GANs. *Neurocomputing*, 2020.

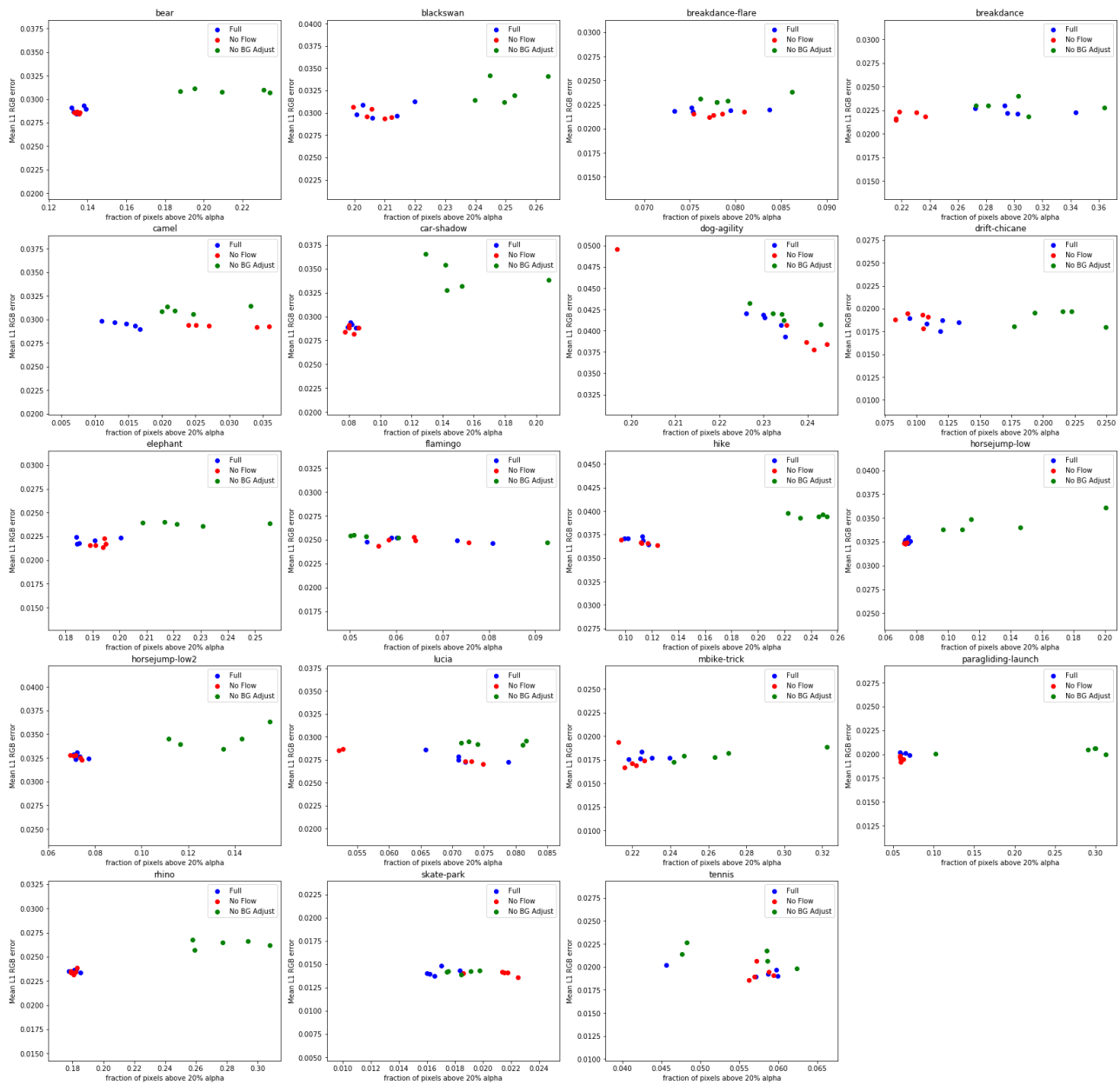


Figure 1. **Reconstruction error vs. matte sparsity for multiple runs on DAVIS.** Each graph measures the L1 reconstruction error against the fraction of pixels above 20% alpha. Lower is better for both axes. A pixel above 20% alpha is considered “visible,” so this value gives a measurement of amount of clutter in the matte. Each video is run with 5 random seeds for each of 3 ablation conditions: full method, no flow input or flow loss (“No Flow”), and no background offset or brightness adjustment (“No BG Adjust”). Note that “No BG Adjust” tends to produce more matte clutter, especially for videos with considerable camera motion (e.g. “horsejump-low”, “hike”).

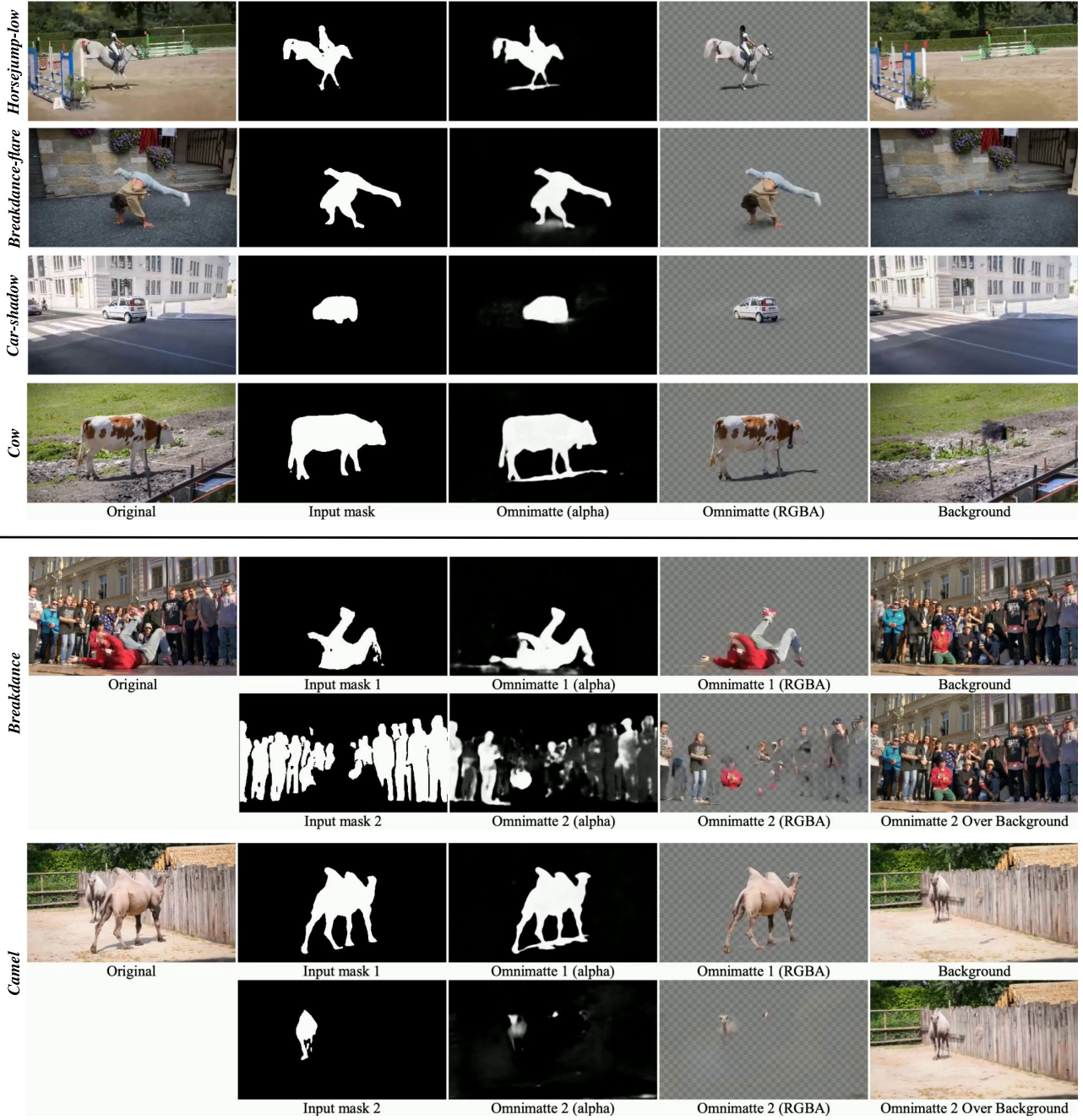


Figure 2. **More Omnimatte results on DAVIS [7].** We show, from left to right, the original RGB video frame, the input mask, the predicted omnimatte’s alpha and RGBA visualizations, and the predicted background layer. Our method successfully detects shadows in each example, and can also handle inaccurate input masks (“Breakdance”, row 1: arm missing in input mask is included in prediction). In “Camel” row 2, the camel in the back is largely static, which results in it being primarily captured by the background layer.

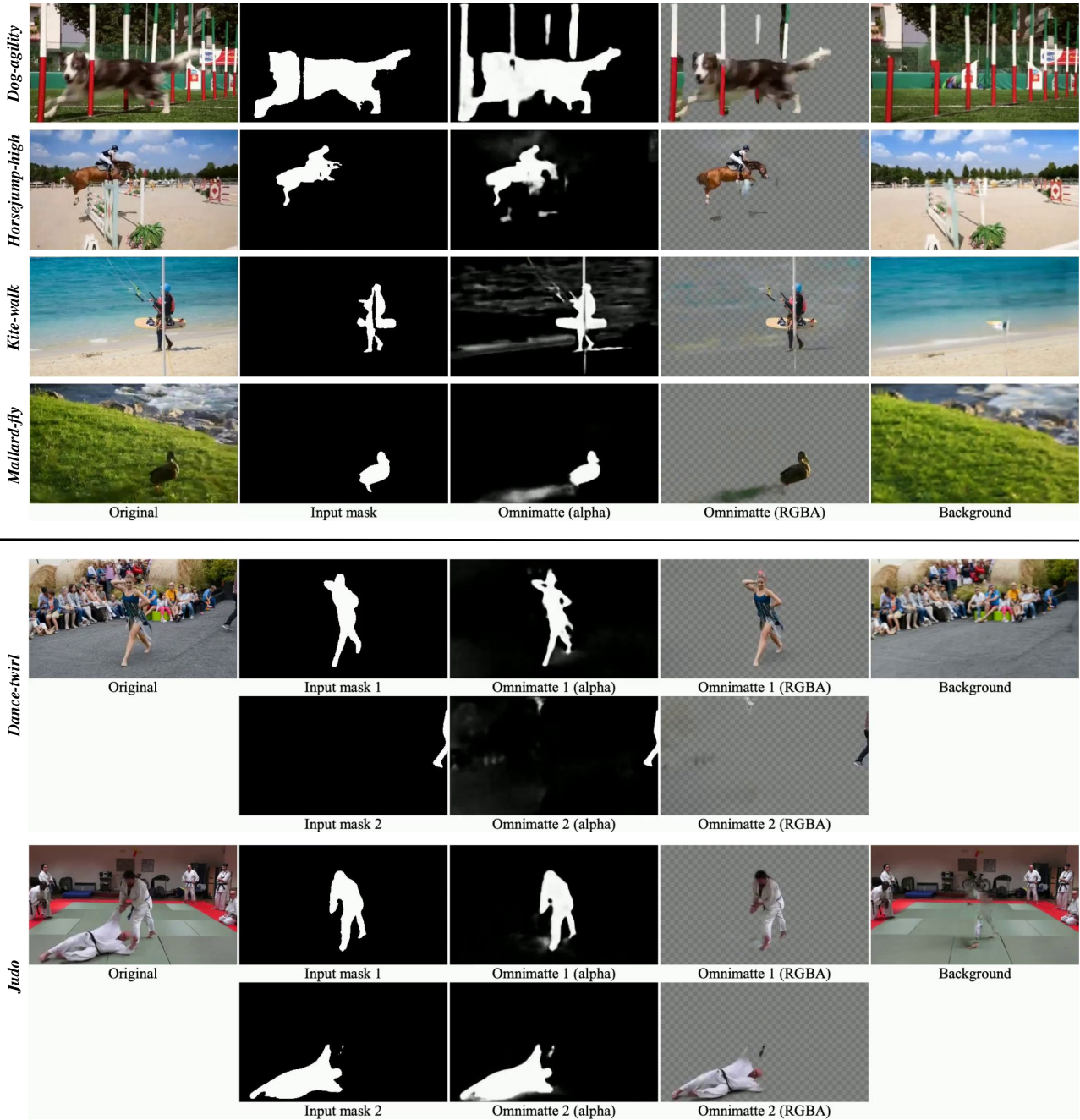


Figure 3. **More Omnimatte results on DAVIS [7].** We show, from left to right, the original RGB video frame, the input mask, the predicted omnimatte’s alpha and RGBA visualizations, and the predicted background layer. We demonstrate our method on a variety of challenging scenes and effects, including shaking poles (“Dog-agility”) and detached shadows (“Horsejump-high”). Our method manages to capture loose clothing not present in the input mask (“Dance-twirl”). The ocean waves in “Kite-walk” are reconstructed in the foreground object layer because they cannot be represented by the static background layer.



Figure 4. **More Omnimatte results on DAVIS [7].** We show, from left to right, the original RGB video frame, the input mask, the predicted omnimatte’s alpha and RGBA visualizations, and the predicted background layer. Our method is able to capture large related objects or shadows (“Paragliding”, “Paragliding-launch”). The water ripples in “Mallard-water” are reconstructed in the foreground object layer because they cannot be represented by the static background layer.

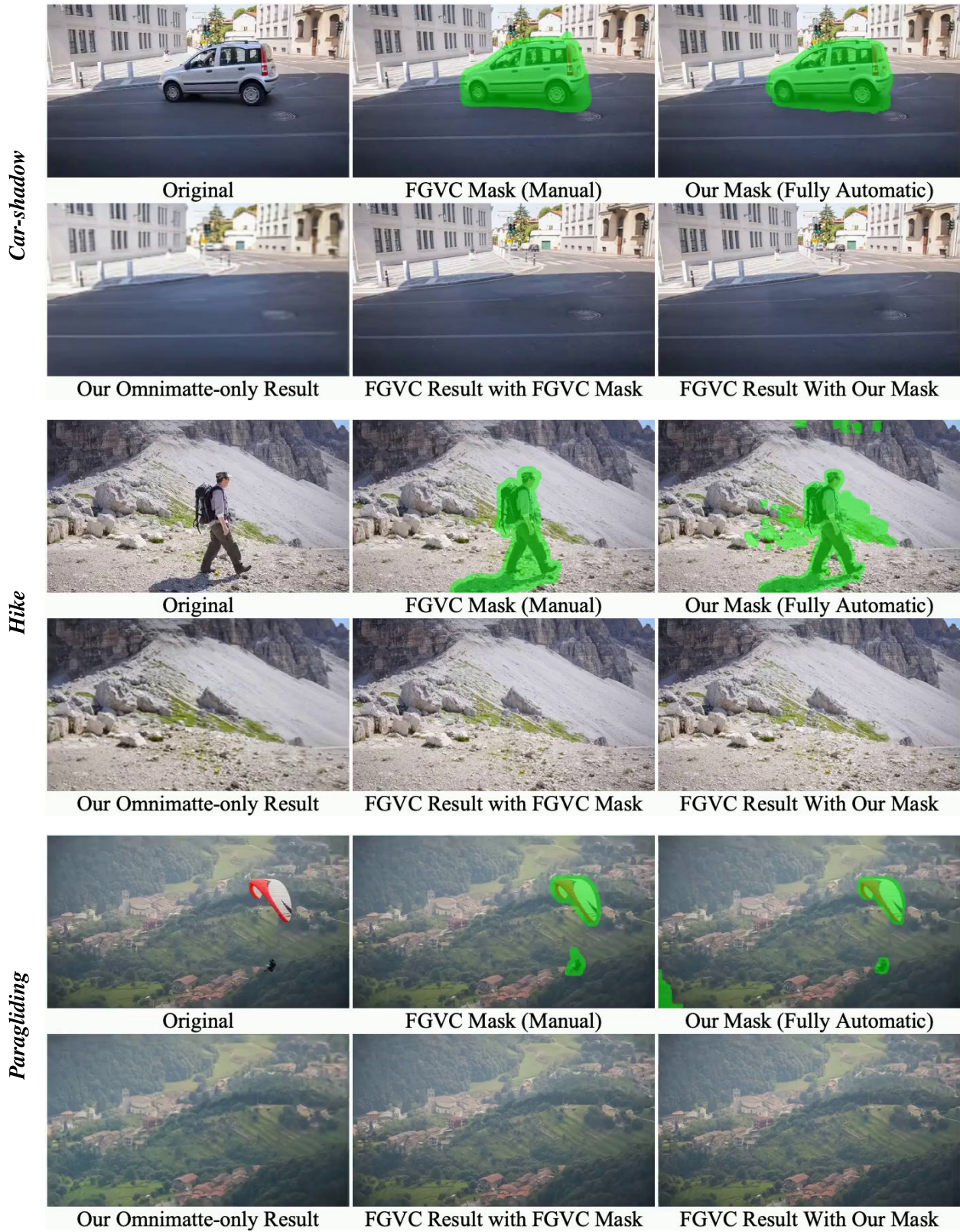


Figure 5. **More object removal comparisons with FGVC [2].** For each example, we show in the top row from left to right: original RGB video frame, the input mask used for FGVC which was manually generated, and our input mask which was automatically generated by thresholding the predicted Omnimatte; in the bottom row: our object removal result by excluding the object layer from composition, the FGVC result using their manual input mask, and the FGVC result using our automatic input mask. We achieve comparable results using our automatic method.

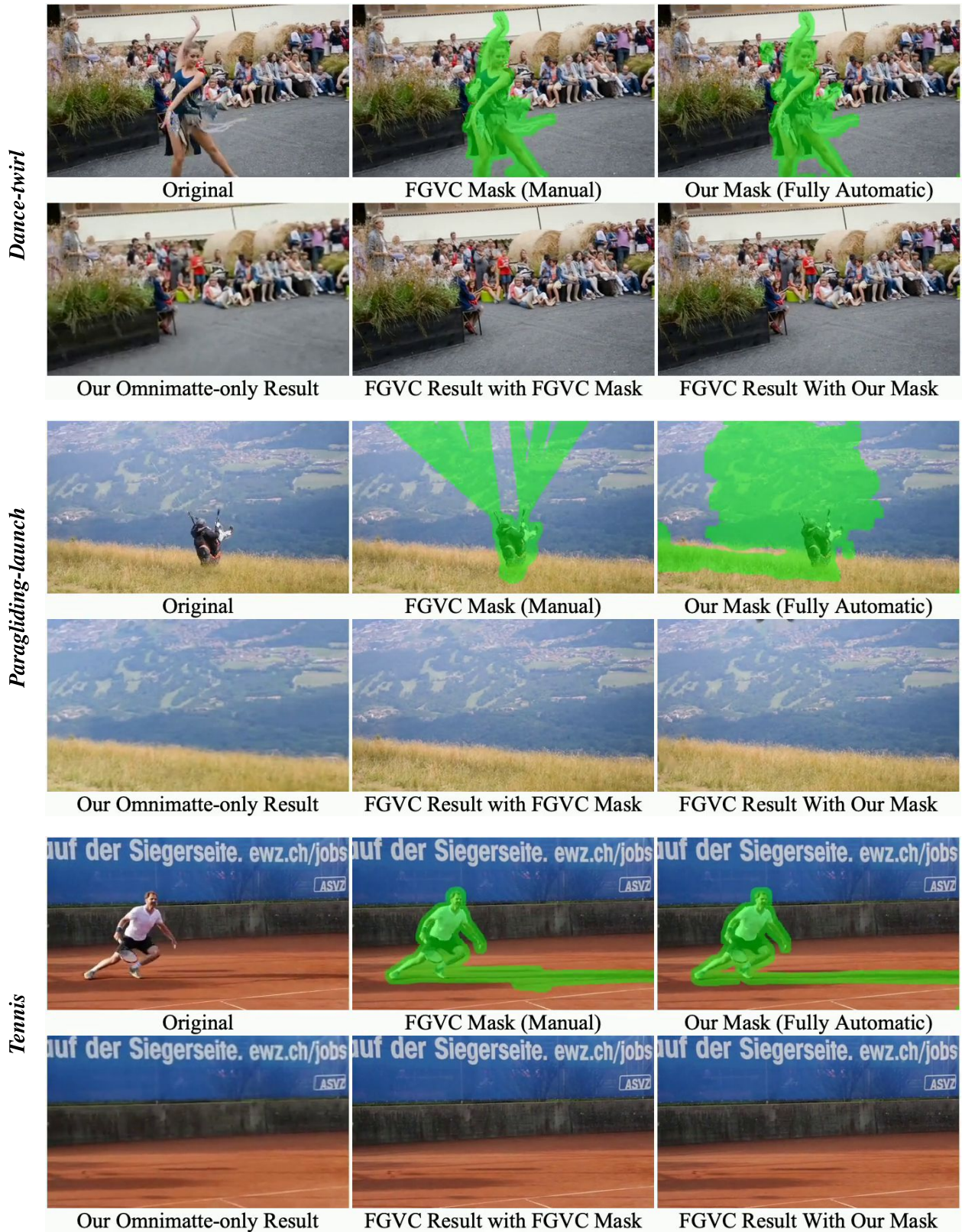


Figure 6. **More object removal comparisons with FGVC [2].** For each example, we show in the top row from left to right: original RGB video frame, the input mask used for FGVC which was manually generated, and our input mask which was automatically generated by thresholding the predicted Omnimatte; in the bottom row: our object removal result by excluding the object layer from composition, the FGVC result using their manual input mask, and the FGVC result using our automatic input mask. We achieve comparable results using our automatic method.

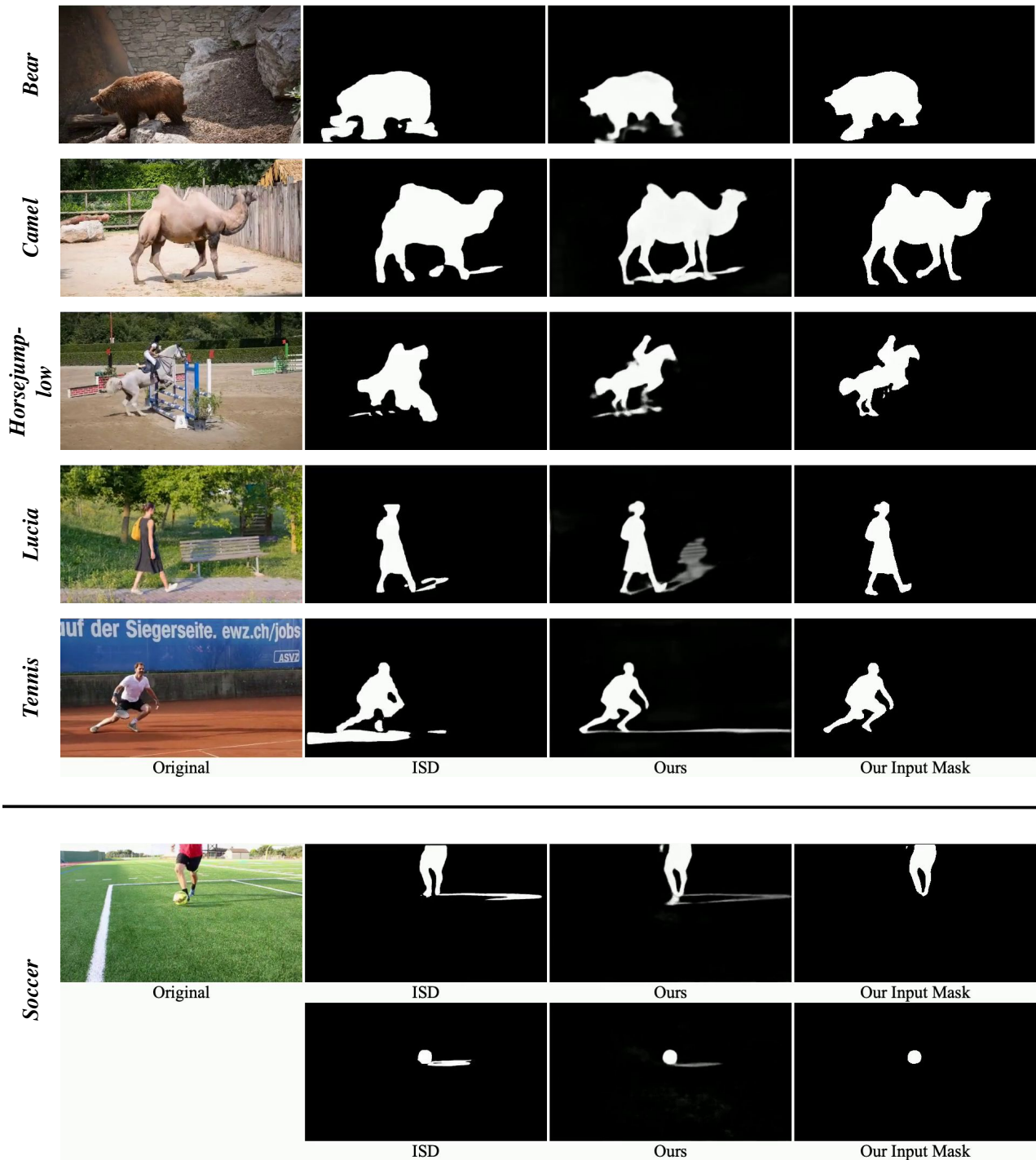


Figure 7. **More shadow detection comparisons with ISD [8]**. For each example, we show from left to right: the original RGB video frame, the object+shadow prediction from ISD, our predicted Omnimatte alpha channel, and our input mask. Our method is able to take advantage of motion in the video, which often enables it to detect shadows more accurately than the single-image, data-driven method, ISD (e.g. in “Tennis”, ISD mistakes other shadows on the ground for the person’s shadow, but our method detects the correct shadow after observing the motion in the scene).

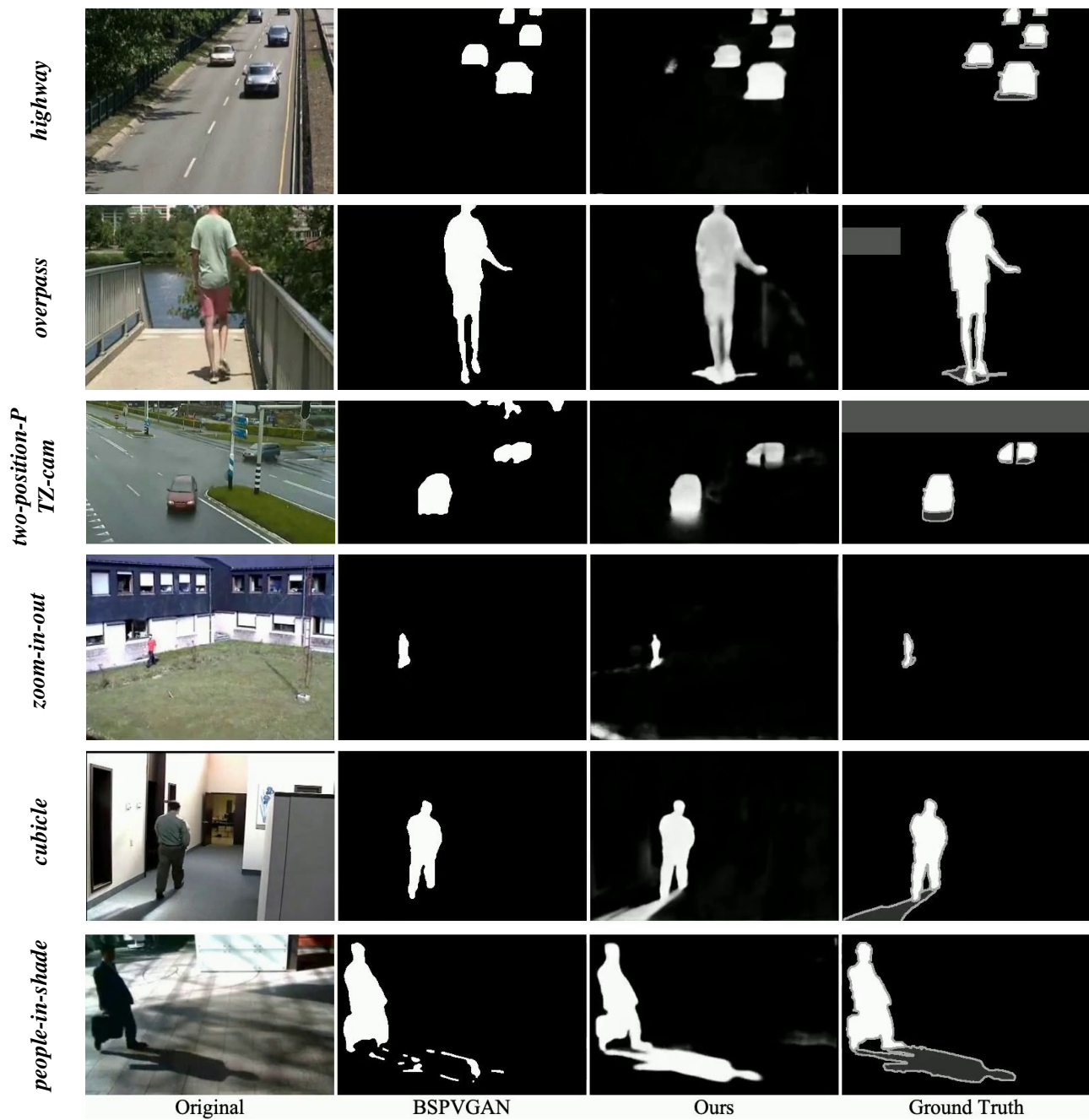


Figure 8. **More change detection comparisons with BSPVGAN [10].** We show comparisons on our subset of CDW-2014 [9]. For each example, we show from left to right: the original RGB video frame, the prediction from BSPVGAN, our predicted Omnimatte alpha channel, and the ground truth mask. The ground truth label definitions are as follows: white pixels = moving objects, dark gray pixels = shadows, light gray pixels = ‘unknown’ (typically at boundaries or unannotated regions). Our method is more successful at detecting shadows.

5

MAST: A Memory-Augmented Self-Supervised Tracker

This work was presented at the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020.

MAST: A Memory-Augmented Self-Supervised Tracker

Zihang Lai Erika Lu Weidi Xie
Visual Geometry Group, Department of Engineering Science
University of Oxford
{zlai, erika, weidi}@robots.ox.ac.uk

Abstract

Recent interest in self-supervised dense tracking has yielded rapid progress, but performance still remains far from supervised methods. We propose a dense tracking model trained on videos **without any annotations** that surpasses previous self-supervised methods on existing benchmarks by a significant margin (+15%), and achieves performance comparable to supervised methods. In this paper, we first reassess the traditional choices used for self-supervised training and reconstruction loss by conducting thorough experiments that finally elucidate the optimal choices. Second, we further improve on existing methods by augmenting our architecture with a crucial memory component. Third, we benchmark on large-scale semi-supervised video object segmentation (aka. dense tracking), and propose a new metric: generalizability. Our first two contributions yield a self-supervised network that for the **first time** is competitive with supervised methods on standard evaluation metrics of dense tracking. When measuring generalizability, we show self-supervised approaches are actually superior to the majority of supervised methods. We believe this new generalizability metric can better capture the real-world use-cases for dense tracking, and will spur new interest in this research direction. Our code will be released at <https://github.com/zlai0/MAST>.

1. Introduction

Although the working mechanisms of the human visual system remain somewhat obscure at the level of neurophysiology, it is a consensus that tracking objects is a fundamental ability that a baby starts developing at two to three months of age [5, 34, 58]. Similarly, in computer vision systems, tracking plays key roles in many applications ranging from autonomous driving to video surveillance.

Given arbitrary objects defined in the first frame, a tracking algorithm aims to relocate the same object throughout the entire video sequence. In the literature, tracking can be cast into two categories: the first is Visual Object Tracking (VOT) [35], where the goal is to relocalize objects

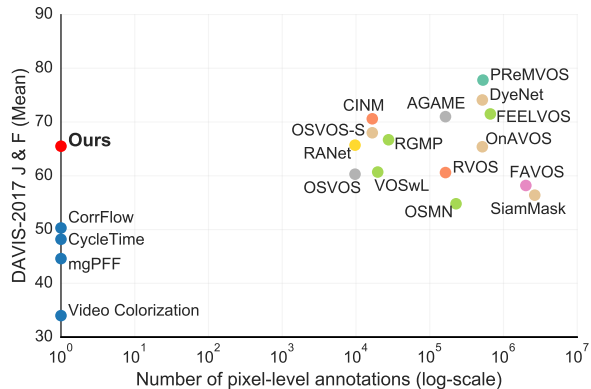


Figure 1: Comparison with other recent works on the DAVIS-2017 benchmarks, i.e. dense tracking or semi-supervised video segmentation given the first frame annotation. The proposed approach significantly outperforms the other self-supervised approaches, and even comparable to approaches trained with heavy supervision on ImageNet, COCO, Pascal, DAVIS, Youtube-VOS. In the x-axis, we only count pixel-wise segmentation.

Notation: CINM [3], OSVOS [6], AGAME [28], VOSwL [31], FAVOS [8], mgPFF [33], CorrFlow [37], DyeNet [39], PReMVOS [41], OSVOS-S [42], RGMP [44], RVOS [54], FEELVOS [56], OnAVOS [57], Video Colorization [59], SiamMask [61], CycleTime [64], RANet [65], OSMN [73].

with bounding boxes throughout the video; the other aims for more fine-grained tracking, i.e. relocalize the objects with pixel-level segmentation masks, also known as Semi-supervised Video Object Segmentation (Semi-VOS) [48]. In this paper, we focus on the latter case, and will refer to it interchangeably with *dense tracking* from here on.

In order to train such dense tracking systems, most recent approaches rely on supervised training with extensive human annotations (see Figure 1). For instance, an ImageNet [10] pre-trained ResNet [18] is typically adopted as a feature encoder, and further fine-tuned on images or video frames annotated with fine-grained, pixelwise segmentation masks, e.g. COCO [40], Pascal [13], DAVIS [48] and YouTube-VOS [71]. Despite their success, this top-down training scheme seems counter-intuitive when considering the development of the human visual system, as infants can track and follow slow-moving objects before they are able to map objects to semantic meanings. With this evidence, it

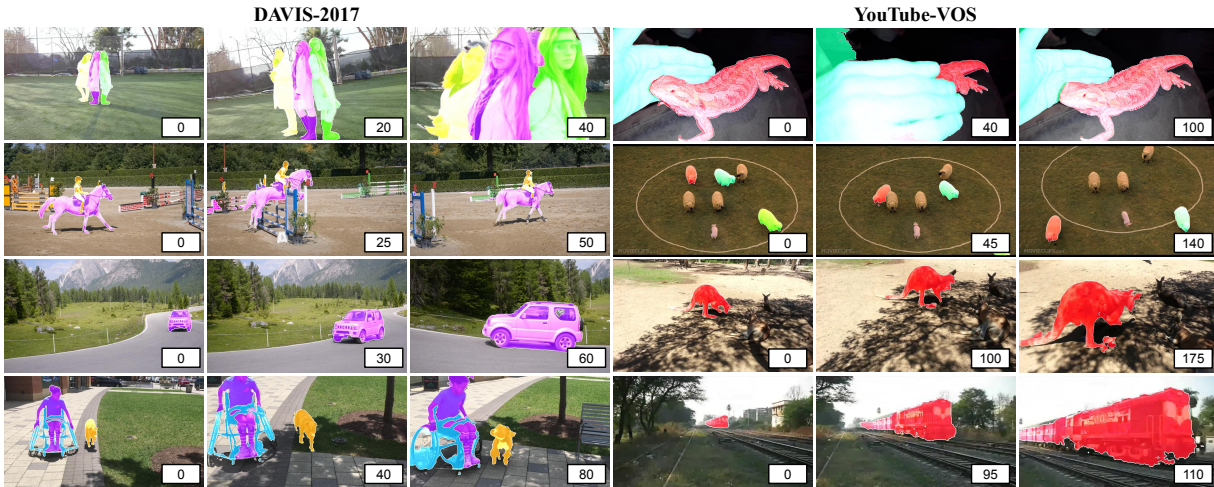


Figure 2: **Train once, test on multiple datasets:** Qualitative results from our *self-supervised dense tracking model* on DAVIS-2017 and YouTube-VOS dataset. The number on the top left refers to the frame number in the video. For all examples, the mask of the 0th frame is given, and the task is to track the objects along with the video. Our self-supervised tracking model is able to deal with challenging scenarios, such as large camera motion, occlusion and disocclusion, large deformation and scale variation.

is unlikely the case that humans develop their tracking ability in a top-down manner (supervised by semantics), at least not at the early-stage development of the visual system.

In contrast to the aforementioned approaches based on heavy supervision, self-supervised methods [37, 59, 60, 64] have recently been introduced, leading to more neurophysiologically intuitive directions. While not requiring any labeled data, the performance of these methods is still far from that of supervised methods (Figure 1).

We continue in the vein of self-supervised methods and propose an improved tracker, which we call *Memory-Augmented Self-Supervised Tracker* (MAST). Similar to previous self-supervised methods, our model performs tracking by learning a feature representation that enables robust pixel-wise correspondences between frames; it then propagates a given segmentation mask to subsequent frames based on the correspondences. We make three main contributions: *first*, we reassess the traditional choices used for self-supervised training and reconstruction loss by conducting thorough experiments to finally determine the optimal choices. *Second*, to resolve the challenge of tracker drift (*i.e.* as the object changes appearance or becomes occluded, each subsequent prediction becomes less accurate if propagated only from recent frames), we further improve on existing methods by augmenting our architecture with a crucial memory component. We design a coarse-to-fine approach that is necessary to efficiently access the memory bank: a two-step attention mechanism first coarsely searches for candidate windows, and then computes fine-grained matching. We conduct experiments to analyze our choice of memory frames, showing that both short- and long-term memory are crucial for good performance. *Third*, we benchmark on large-scale video segmentation datasets and propose a new metric, *i.e.* generalizability, with the goal

of measuring the performance gap between tracking seen and unseen categories, which we believe better captures the real-world use-cases for category-agnostic tracking.

The result of the first two contributions is a self-supervised network that surpasses all existing approaches by a significant margin on DAVIS-2017 (15%) and YouTube-VOS (17%) benchmarks, making it competitive with supervised methods *for the first time*. Our results show that a strong representation for tracking can be learned without using any semantic annotations, echoing the early-stage development of the human visual system. Beyond significantly narrowing the gap with supervised methods on the existing metrics, we also demonstrate the *superiority* of self-supervised approaches over supervised methods on generalizability. On the unseen categories of YouTube-VOS benchmark, we surpass PreMVOS [41], the 2018 challenge winner algorithm trained on massive segmentation datasets. Furthermore, when we analyze the drop in performance between seen and unseen categories, we show that our method (along with other self-supervised methods) has a significantly smaller *generalization gap* than supervised methods. These results show that contrary to the popular belief that self-supervised methods are not yet useful due to their weaker performance, their greater generalization capability (due to not being at risk of overfitting to labels) is actually a more desirable quality when being deployed in real-world settings, where the domain gap can be significant.

2. Related Work

Dense tracking (*aka. semi-supervised video segmentation*) has typically been approached in one of two ways: propagation-based or detection/segmentation-based. The

former approaches formulate the dense tracking task as a mask propagation problem from the first frame to the consecutive frames. To leverage the temporal consistency between two adjacent frames, many propagation-based methods often try to establish dense correspondences with optical flow or metric learning [20, 21, 29, 41, 56]. However, computing optical flow remains a challenging, yet unsolved problem. Our method relaxes the constraint of optical flow’s one-to-one brightness constancy constraint and spatial smoothness, allowing each query pixel to potentially build correspondence with multiple reference pixels. On the other hand, detection/segmentation-based approaches address the tracking task with sophisticated detection or segmentation networks, but since these models are usually not class-agnostic during training, they often have to be fine-tuned on the first frame of the target video during inference [6, 41, 42], whereas our method requires no fine-tuning.

Self-supervised learning on videos has generated fruitful research in recent years. Due to the abundance of online data [1, 4, 11, 14, 15, 22, 24, 25, 26, 27, 32, 38, 43, 59, 63, 67, 68], various ideas have been explored to learn representations by exploiting the spatio-temporal information in videos. [14, 43, 66] exploit spatio-temporal ordering for learning video representations. Recently, Han *et al.* [17] learn strong video representations for action recognition by self-supervised contrastive learning on raw videos. Of more relevance, [37, 59] have recently leveraged the natural temporal coherency of color in videos, to train a network for tracking and correspondence related tasks. We discuss these works in more detail in Section 3.1. In this work, we propose to augment the self-supervised tracking algorithms with a differentiable memory module. We also rectify some flaws in their training process.

Memory-augmented models refer to the computational architecture that has access to a memory repository for prediction. Such models typically involve an internal memory implicitly updated in a recurrent process, *e.g.* LSTM [19] and GRU [9], or an explicit memory that can be read or written with an attention-based procedure [2, 12, 16, 36, 51, 53, 62, 70]. Memory models have been used for many applications, including reading comprehension [51], summarization [50], tracking [69], video understanding [7], and image and video captioning [70, 74]. In dense visual tracking, the popular memory-augmented models treat key frames as memory [45], and use attention mechanisms to read from the memory. Despite its effectiveness, the process of computing attention either does not scale to multiple frames or is unable to process high-resolution frames, due to the computational bottleneck in hardware, *e.g.* physical memory. In this work, we propose a scalable way to process high-resolution information in a coarse-to-fine manner. The model enables dynamic localization of salient regions, and fine-grained processing is only required for a small fraction

of the memory bank.

3. Method

The proposed dense tracking system, MAST (Memory-Augmented Self-Supervised Tracker), is a conceptually simple model for dense tracking that can be trained with self-supervised learning, *i.e.* **zero manual annotation** is required during training, and an object mask is only required for the first frame during inference. In Section 3.1, we provide relevant background of previous self-supervised dense tracking algorithms, and terminologies that will be used in later sections. Next, in Section 3.2, we pinpoint weaknesses in these works and propose improvements to the training signals. Finally, in Section 3.3, we propose memory augmentation as an extension to existing self-supervised trackers.

3.1. Background

In this section, we review previous papers that are closely related to this work [37, 59]. In general, the goal of self-supervised tracking is to learn feature representations that enable robust correspondence matching. During training, a proxy task is posed as reconstructing a target frame (I_t) by linearly combining pixels from a reference frame (I_{t-1}), with the weights measuring the strength of correspondence between pixels.

Specifically, a triplet ($\{Q_t, K_t, V_t\}$) exists for each input frame I_t , referring to *Query*, *Key*, and *Value*. In order to reconstruct a pixel i in the t -th frame (\hat{I}_t^i), an *Attention* mechanism is used for *copying* pixels from a subset of previous frames in the original sequence. This procedure is formalized as:

$$\hat{I}_t^i = \sum_j A_t^{ij} V_{t-1}^j \quad (1)$$

$$A_t^{ij} = \frac{\exp\langle Q_t^i, K_{t-1}^j \rangle}{\sum_p \exp\langle Q_t^i, K_{t-1}^p \rangle} \quad (2)$$

where $\langle \cdot, \cdot \rangle$ refers to the dot product between two vectors, query (Q) and key (K) are feature representations computed by passing the target frame I_t to a Siamese ConvNet $\Phi(\cdot; \theta)$, *i.e.* $Q_t = K_t = \Phi(I_t; \theta)$, A_t is the affinity matrix representing the feature similarity between pixel I_t^i and I_{t-1}^j , value (V) is the raw reference frame (I_{t-1}) during the training stage, and instance segmentation mask during inference, achieving reconstruction or dense tracking respectively.

A key element in self-supervised learning is to set the proper *information bottleneck*, or the choice of what input information to withhold for learning the desired feature representation and avoiding trivial solutions. For example, in the reconstruction-by-copying task, an obvious shortcut is that the pixel in I_t can learn to match any pixel in I_{t-1} with the *exact same color*, yet not necessarily correspond

to the same object. To circumvent such learning shortcuts, Vondrick *et al.* [59] intentionally drop the color information from the input frames. Lai and Xie [37] further show that a simple channel dropout can be more effective.

3.2. Improved Reconstruction Objective

In this section, we reassess the choices made in previous self-supervised dense tracking works and provide intuition for our optimal choices, which we empirically support in Section 5.

3.2.1 Decorrelated Color Space

Extensive experiments in the human visual system have shown that colors can be seen as combinations of the primary colors, namely red (R), green (G) and blue (B). For this reason, most of the cameras and emissive color displays represent pixels as a triplet of intensities: $(R, G, B) \in \mathcal{R}^3$. However, a disadvantage of the RGB representation is that the channels tend to be extremely correlated [49], as shown in Figure 3. In this case, the channel dropout proposed in [37] is unlikely to behave as an effective information bottleneck, since the dropped channel can almost always be determined by one of the remaining channels.

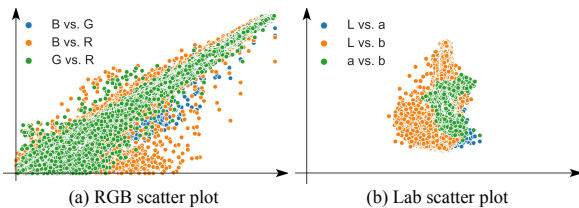


Figure 3: Correlation between channels of RGB and *Lab* colorspace. We randomly take 100,000 pixels from 65 frames in a sequence (snowboard) in the DAVIS dataset and plot the relative relationships between RGB channels. This phenomena generally holds for all natural images [49], due to the fact that all of the channels include a representation of brightness. Values are normalized for visualization purposes.

To remedy this limitation, we hypothesize that dropout in the decorrelated representations (*e.g.* *Lab*) would force the model to learn invariances suitable for self-supervised dense tracking; *i.e.* if the model cannot predict the missing channel from the observed channels, it is forced to learn a more robust representation rather than relying on local color information.

3.2.2 Classification vs. Regression

In the recent literature on colorization and generative models [46, 75], colors were quantized into discrete classes and treated as a multinomial distribution, since generating images or predicting colors from grayscale images is usually a non-deterministic problem; *e.g.* the color of a car can reasonably be red or white. However, this convention is suboptimal for self-supervised learning of correspondences,

as we are not trying to generate colors for each pixel, but rather, estimate a precise relocation of pixels in the reference frames. More importantly, quantizing the colors leads to an information loss that can be crucial for learning high-quality correspondences.

We conjecture that directly optimizing a regression loss between the reconstructed frame (\hat{I}_t) and real frame (I_t) will provide more discriminative training signals. In this work, the objective \mathcal{L} is defined as the Huber Loss:

$$\mathcal{L} = \frac{1}{n} \sum_i z_i \quad (3)$$

where

$$z_i = \begin{cases} 0.5(\hat{\mathbf{I}}_t^i - \mathbf{I}_t^i)^2, & \text{if } |\hat{\mathbf{I}}_t^i - \mathbf{I}_t^i| < 1 \\ |\hat{\mathbf{I}}_t^i - \mathbf{I}_t^i| - 0.5, & \text{otherwise} \end{cases} \quad (4)$$

where $\hat{\mathbf{I}}_t^i \in \mathcal{R}^3$ refers to RGB or *Lab*, normalized to the range [-1,1] in the reconstructed frame that is copied from pixels in the reference frame I_{t-1} , and I_t is the real frame at time point t .

3.3. Memory-Augmented Tracking

So far we have discussed the straightforward attention-based mechanism for propagating a mask from a single previous frame. However, as predictions are made recursively, errors caused by object occlusion and disocclusion tend to accumulate and eventually degrade the subsequent predictions. To resolve this issue, we propose an attention-based tracker that efficiently makes use of *multiple* reference frames.

3.3.1 Multi-frame tracker

An overview of our tracking model is shown in Figure 4. To summarize the tracking process: given the present frame and multiple past frames (memory bank) as input, we first compute the query (Q) for the present frame and keys (K) for all frames in memory. Here, we follow the general procedure in previous works as described in Section 3.1, where K and Q are computed from a shared-weight feature extractor and V is equal to the input frame (during training) or object mask (during testing). The computed affinity between Q and all the keys (K) in memory is then used to make a prediction for each query pixel depending on V. Note we don't put any weights on the reference frames, as this should be encoded in the affinity matrix (*e.g.* when a target and reference frame are dis-similar, the corresponding similarity value will be naturally low; thus the reference label will have less contribution to the labeling of a target pixel).

The decision of which pixels to include in K is crucial for good performance. Including all pixels previously seen is far too computationally expensive due to the quadratic explosion of the affinity matrix (*e.g.* the network of [37] produces affinity matrices with more than 1 billion elements

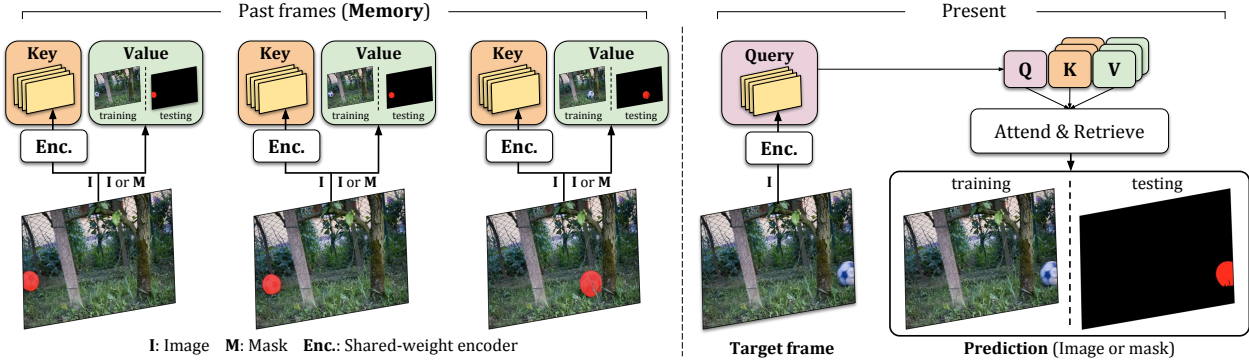


Figure 4: Structure of MAST. The current frame is used to compute **query** to attend and retrieve from memory (**key** & **value**). During training, we use raw video frame as **value** for self-supervision. Once the encoder is trained, we use instance mask as **value**. See Section 3.3 for details.

for 480p videos). To reduce computation, [37] exploit temporal smoothness in videos and apply restricted attention, only computing the affinity with pixels in a ROI around the query pixel location. However, the temporal smoothness assumption holds only for temporally close frames.

To efficiently process temporally distant frames, we propose a two-step attention mechanism. The first stage involves coarse pixel-matching with the frames in the memory bank to determine which ROIs are likely to contain good matches with the query pixel. In the second stage, we extract the ROIs and compute fine-grained pixel matching, as described in Section 3.1. Overall, the process can be summarized in Algorithm 1.

Algorithm 1 MAST

- 1: Choose m reference frames Q_1, Q_2, \dots, Q_m
- 2: Localize ROI R_1, R_2, \dots, R_m according to 3.3.2 (Eq. 5 and 6) for each of the reference frames
- 3: Compute similarity matrix $A_t^{ij} = \langle Q^j, R_t^i \rangle$ between target frame Q and each ROI.
- 4: Output: pixel's label is determined by aggregating the labels of the ROI pixels (weighted by its affinity score).

3.3.2 ROI Localization

The goal of ROI localization is to estimate the candidate windows non-locally from memory banks. As shown in Figure 5, this can be achieved through restricted attention with varying dilation rate.

Intuitively, for short-term memory (temporally close frames), dilation is not required since spatial-temporal coherence naturally exists in videos; thus ROI localization becomes restricted attention (similar to [37]). However, for long-term memory, we aim to account for the fact that objects can potentially appear anywhere in the reference frames. We unify both scenarios into a single framework for learning ROI localization.

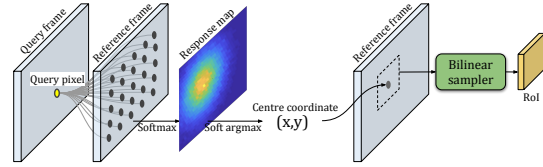


Figure 5: ROI Localization.

Formally, for the query pixel i in I_t , to localize the ROI from frame (I_{t-N}), we first compute in parallel $H_{t-N,x,y}^i$, the similarity heatmap between i and all candidate pixels in the dilated window:

$$H_{t-N,x,y}^i = \text{softmax}(Q_t^i \cdot \text{im2col}(K_{t-N}^i, \gamma_{t-N})) \quad (5)$$

where γ_{t-N} refers to the dilation rate for window sampling in frame I_{t-N} , and im2col refers to an operation that transforms the input feature map into a matrix based on dilation rate. Specifically, in our experiments, the dilation rate is proportional to the temporal distance between the present frame and the past frames in the memory bank, i.e. $\gamma_{t-N} \propto N$. We use $\gamma_{t-N} = \lceil (t-N)/15 \rceil$.

The center coordinates for ROIs can be then computed via a *soft-argmax* operation:

$$P_{x,y}^i = \sum_{x,y} H_{x,y}^i * C \quad (6)$$

where $P_{x,y}^i$ is the estimated center location of the candidate window in frame I_{t-N} for query pixel I_t^i , and C refers to the grid coordinates (x, y) corresponding to the pixels in the window from im2col . The resampled Key (\hat{K}_{t-N}^i) for pixel I_t^i can be extracted with a bilinear sampler [23]. With all the candidate Keys dynamically sampled from different reference frames of the memory bank, we compute fine-grained matching scores only with these localized Keys, resembling a restricted attention in a non-local manner. Therefore, with the proposed design, the memory-augmented model can efficiently access high-resolution information for correspondence matching, without incurring large physical memory costs.

4. Implementation Details

Training: For fair comparison, we adopt as our feature encoder the same architecture (ResNet18) as [37] in all experiments (as shown in Supplementary Material). The network produces feature embeddings with a spatial resolution 1/4 of the original image. The model is trained in a completely self-supervised manner, meaning the model is initialized with random weights, and we do *not* use any information other than raw video sequences. We report main results on two training datasets: OxUvA [52] and YouTube-VOS (both raw videos only). We report the first for fair comparison with the state-of-the-art method [37] and the second for maximum performance. As pre-processing, we resize all frames to $256 \times 256 \times 3$. In all of our experiments, we use I_0, I_5 (only if the index for the current frame is larger than 5) as long term memory, and $I_{t-5}, I_{t-3}, I_{t-1}$ as short term memory. Empirically, we find the choice of frame number has small impact on performance, but using both long and short term memory is essential (See appendix).

During training, we first pretrain the network with a pair of input frames, *i.e.* one reference frame and one target frame are fed as inputs. One of the color channels is randomly dropped with probability $p = 0.5$. We train our model end-to-end using a batch size of 24 for 1M iterations with the Adam optimizer. The initial learning rate is set to $1e-3$, and halved after 0.4M, 0.6M and 0.8M iterations. We then finetune the model using multiple reference frames (our full memory-augmented model) with a small learning rate of $2e-5$ for another 1M iterations. As discussed in Section 3.2.2, the model is trained with a photometric loss between the reconstruction and the true frame.

Inference: We use the trained feature encoder to compute the affinity matrix between pixels in the target frame and those in the reference frames. The affinity matrix is then used to propagate the desired pixel-level entities, such as instance masks in the dense tracking case, as described in Algorithm 1.

Image Feature Alignment: Due to memory constraints, the supervision signals in previous methods were all defined on bilinearly downsampled images. As shown in Figure 6(a), this introduces a misalignment between strided convolution layers and images from naïve bilinear downsampling. We handle this spatial misalignment between feature embedding and image by directly sampling at the strided convolution centers. This seemingly minor change actually brings significant improvement to the downstream tracking task (Table 4).

5. Experiments

We benchmark our model on two public benchmarks: DAVIS-2017 [48] and the current largest video segmentation dataset, YouTube-VOS [71]. The former contains 150

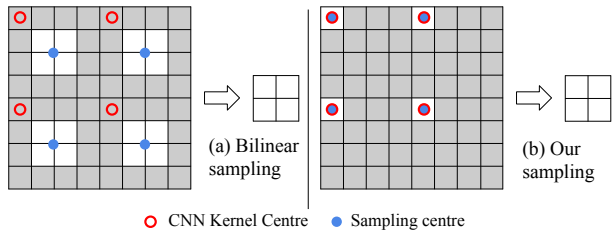


Figure 6: Image Feature Alignment. We fix the misalignment between strided convolution layers (with kernel centered at red circle) and images from naïve bilinear downsampling (with kernel centered at blue dot) by sampling directly at the strided convolution centers.

HD videos with over 30K manual instance segmentations, and the latter has over 4000 HD videos of 90 semantic categories, totalling over 190k instance segmentations. For both datasets, we benchmark the proposed self-supervised learning architecture (MAST) on the official semi-supervised video segmentation setting (*aka.* dense tracking), where a ground truth instance segmentation mask is given for the first frame, and the objective is to propagate the mask to subsequent frames. In Section 5.1, we report performance of our full model and several ablated models on the DAVIS benchmark. Next, in Section 5.2, we analyze the *generalizability* of our model by benchmarking on the large-scale YouTube-VOS dataset.

Standard evaluation metrics. We use region similarity (\mathcal{J}) and contour accuracy (\mathcal{F}) to evaluate the tracked instance masks [47].

Generalizability metrics To demonstrate the generalizability of tracking algorithms in category-agnostic scenarios, *i.e.* the categories in training set and testing set are disjoint, YouTube-VOS also explicitly benchmarks the performances on unseen categories. We therefore evaluate a *generalization gap* in Section 5.2.1, which is defined as the average performance difference between seen and unseen object classes:

$$\text{Gen. Gap} = \frac{(\mathcal{J}_{\text{seen}} - \mathcal{J}_{\text{unseen}}) + (\mathcal{F}_{\text{seen}} - \mathcal{F}_{\text{unseen}})}{2} \quad (7)$$

Note, the proposed metric aims to explicitly penalize the case where the performance on seen outperforms unseen by large margins, while at the same time provide a reward when the performance on unseen categories is higher than on seen ones.

5.1. Video Segmentation on DAVIS-2017

5.1.1 Main results

In Table 1, we compare MAST with previous approaches on the DAVIS-2017 benchmark. Two phenomena can be observed: *first*, our proposed model clearly outperforms all other self-supervised methods, surpassing previous state-of-the-art CorrFlow by a significant margin (65.5 vs 50.3 on

Method	Backbone	Supervised	Dataset (Size)	$\mathcal{J}\&\mathcal{F}$ (Mean) \uparrow	\mathcal{J} (Mean) \uparrow	\mathcal{J} (Recall) \uparrow	\mathcal{F} (Mean) \uparrow	\mathcal{F} (Recall) \uparrow
Vid. Color. [59]	ResNet-18	\times	Kinetics (800 hours)	34.0	34.6	34.1	32.7	26.8
CycleTime [†] [64]	ResNet-50	\times	VLOG (344 hours)	48.7	46.4	50.0	50.0	48.0
CorrFlow [†] [37]	ResNet-18	\times	OxUvA (14 hours)	50.3	48.4	53.2	52.2	56.0
UVC* [72]	ResNet-18	\times	Kinetics (800 hours)	59.5	57.7	68.3	61.3	69.8
MAST (Ours)	ResNet-18	\times	OxUvA (14 hours)	63.7	61.2	73.2	66.3	78.3
MAST (Ours)	ResNet-18	\times	YT-VOS (5.58 hours)	65.5	63.3	73.2	67.6	77.7
ImageNet [18]	ResNet-50	\checkmark	I (1.28M, 0)	49.7	50.3	-	49.0	-
OSMN [73]	VGG-16	\checkmark	ICD (1.28M, 227k)	54.8	52.5	60.9	57.1	66.1
SiamMask [61]	ResNet-50	\checkmark	IVCY (1.28M, 2.7M)	56.4	54.3	62.8	58.5	67.5
OSVOS [6]	VGG-16	\checkmark	ID (1.28M, 10k)	60.3	56.6	63.8	63.9	73.8
OnAVOS [57]	ResNet-38	\checkmark	ICPD (1.28M, 517k)	65.4	61.6	67.4	69.1	75.4
OSVOS-S [42]	VGG-16	\checkmark	IPD (1.28M, 17k)	68.0	64.7	74.2	71.3	80.7
FEELVOS [56]	Xception-65	\checkmark	ICDY (1.28M, 663k)	71.5	69.1	79.1	74.0	83.8
PRemVOS [41]	ResNet-101	\checkmark	ICDPM (1.28M, 527k)	77.8	73.9	83.1	81.8	88.9
STM [45]	ResNet-50	\checkmark	IDY (1.28M, 164k)	81.8	79.2	-	84.3	-

Table 1: Video segmentation results on DAVIS-2017 validation set. Dataset notations: I=ImageNet, V = ImageNet-VID, C=COCO, D=DAVIS, M=Mapillary, P=PASCAL-VOC Y=YouTube-VOS. For size of datasets, we report (length of *raw videos*) for self-supervised methods and (#image-level annotations, #pixel-level annotations) for supervised methods. * denotes concurrent work. [†] denotes highest results reported after original publication. Higher values are better.

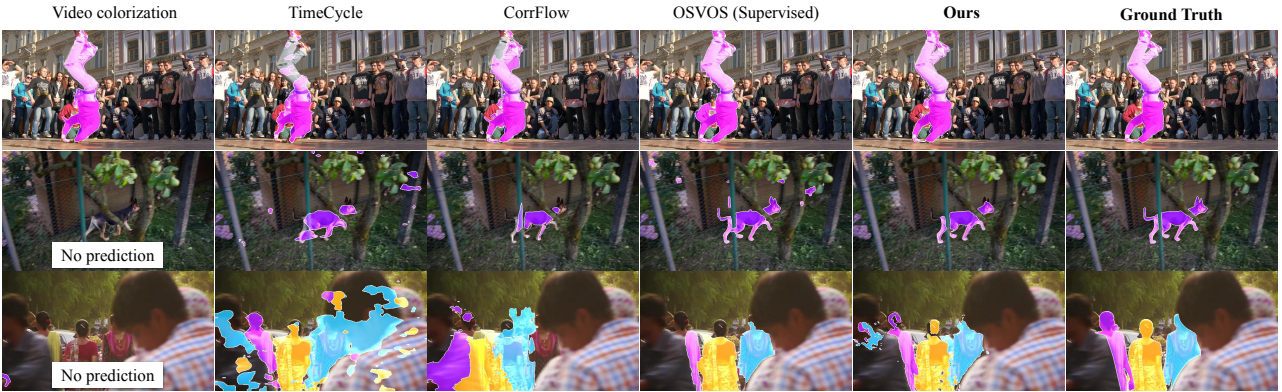


Figure 7: Our method vs. previous self-supervised methods. Other methods show systematic errors in handling occlusions. Row 1: The dancer undergoes large self-occlusion. Row 2: The dog is repeatedly occluded by poles. Row 3: Three women reappear after being occluded by the man in the foreground.

$\mathcal{J}\&\mathcal{F}$). *Second*, despite using only ResNet18 as the feature encoder, our model trained with self-supervised learning can still surpass supervised approaches that use heavier architectures.

5.1.2 Ablation Studies

To examine the effects of different components, we conduct a series of ablation studies by removing one component at a time. All models are trained on OxUvA (except for the analysis on different datasets), and evaluated on DAVIS-2017 semi-supervised video segmentation (*aka.* dense tracking) *without* any finetuning.

Choice of color spaces. As shown in Table 2, we perform different experiments with input frames transformed into different color spaces, *e.g.* RGB, *Lab* or HSV. We find that the MAST model trained with *Lab* color space always outperforms the other color spaces, validating our conjecture that dropout in a decorrelated color space leads to better feature representations for self-supervised dense tracking,

as explained in Section 3.2.1. Additionally, we compare our default setting with a model trained with cross-color space matching task (shown in Table 3). That means to use a different color space for the input and the training objective, *e.g.* input frames are in RGB, and loss function is defined in *Lab* color space. Interestingly, the performance drops significantly, we hypothesis this can attribute to the fact that all RGB channels include a representation of brightness, making it highly correlate to the luminance in *Lab*, therefore acting as a weak information bottleneck.

Loss functions. As a variation of our training procedure, we experiment with different loss functions: cross entropy loss on the quantized colors, and photometric loss with Huber loss. As shown in Table 2, regression with real-valued photometric loss surpasses classification significantly, validating our conjecture that the information loss during color quantization results in inferior representations for self-supervised tracking (as explained in Section 3.2), due to less discriminative training signals.

Colors.	Loss	$\mathcal{J}(\text{Mean})$	$\mathcal{F}(\text{Mean})$
RGB	Cls.	42.5	45.3
	Reg.	52.7	57.1
HSV	Cls.	32.5	35.3
	Reg.	54.3	58.6
Lab	Cls.	47.1	48.9
	Reg.	61.2	66.3

Table 2: **Training colorspace and loss:** Our final model trained with Lab colorspace with regression loss outperforms all other models on dense tracking task. Higher values are better.

Memory	$\mathcal{J}(\text{Mean})$	$\mathcal{F}(\text{Mean})$
Only long	44.6	48.7
Only short	57.3	61.8
Both	61.2	66.3

Table 5: **Memory length:** Removing either long term or short term memory results in a performance drop.

Image feature alignment. To evaluate the alignment module proposed for aligning features with the original image, we compare it to direct bilinear image downsampling used by CorrFlow [37]. The result in Table 4 shows that our approach achieves about 2.2% higher performance.

Dynamic memory by exploiting more frames. We compare our default network with variants that have only short term memory or long term memory. Results are shown in Table 5. While both short term memory and long term memory alone can make reasonable predictions, the combined model achieves the highest performance. The qualitative predictions (Figures 10 and 7) also confirm that the improvements come from reduced tracker drift. For instance, when severe occlusion occurs, our model is able to attend and retrieve high-resolution information from frames that are temporally distant.

5.2. Youtube Video Object Segmentation

We also evaluate the MAST model on the Youtube-VOS validation split (474 videos with 91 object categories). As no other self-supervised methods have been tested on the benchmark, we directly compare our results with supervised methods. As shown in Table 8, our method outperforms the other self-supervised learning approaches by a significant margin (64.2 vs. 46.6), and even achieves comparable performance to many heavily supervised methods.

5.2.1 Generalizability

As another metric for evaluating category-agnostic tracking, the YouTube-VOS dataset conveniently has separate measures for seen and unseen object categories. We can therefore estimate testing performance on out-of-

Input	Loss	$\mathcal{J}(\text{Mean})$	$\mathcal{F}(\text{Mean})$
Lab	RGB	48.2	52.0
RGB	Lab	46.8	49.9
Lab	Lab	61.2	66.3

Table 3: **Cross color space matching vs. single color space:** Cross color space matching shows inferior results compared to single color space.

Propagation	$\mathcal{J}(\text{Mean})$	$\mathcal{F}(\text{Mean})$
Soft	57.0	61.7
Hard	61.2	66.3
	+4.2	+4.6

Table 6: **Soft vs. hard propagation:** *Quantizing* class probability of each pixel (hard propagation) shows large gains over propagating probability distribution (soft propagation).

I-F Align	$\mathcal{J}(\text{Mean})$	$\mathcal{F}(\text{Mean})$
No	59.1	64.0
Yes	61.2	66.3
	+2.1	+2.3

Table 4: **Image-Feature alignment:** Using the improved Image-Feature alignment implementation improves the results. Higher values are better.

Dataset	$\mathcal{J}(\text{Mean})$	$\mathcal{F}(\text{Mean})$
OxUvA	61.2	66.3
ImageNet VID	60.0	63.9
YouTube-VOS (w/o anno.)	63.3	67.6

Table 7: **Training dataset:** All datasets provide reasonable performance, with O and Y slightly superior. We conjecture that our model gains from higher quality videos and larger object classes in these datasets.

Method	Sup.	Overall \uparrow	Seen		Unseen		Gen. Gap \downarrow
			$\mathcal{J} \uparrow$	$\mathcal{F} \uparrow$	$\mathcal{J} \uparrow$	$\mathcal{F} \uparrow$	
Vid. Color.[59] [†]	\times	38.9	43.1	38.6	36.6	37.4	3.9
CorrFlow[37]	\times	46.6	50.6	46.6	43.8	45.6	3.9
MAST (Ours)	\times	64.2	63.9	64.9	60.3	67.7	0.4
OSMN[73]	\checkmark	51.2	60.0	60.1	40.6	44.0	17.75
MSK[30]	\checkmark	53.1	59.9	59.5	45.0	47.9	13.25
RGMP[44]	\checkmark	53.8	59.5	-	45.2	-	14.3
OnAVOS[57]	\checkmark	55.2	60.1	62.7	46.6	51.4	12.4
RVOS[55]	\checkmark	56.8	63.6	67.2	45.5	51.0	17.15
OSVOS[6]	\checkmark	58.8	59.8	60.5	54.2	60.7	2.7
S2S[71]	\checkmark	64.4	71.0	70.0	55.5	61.2	12.15
PreMVOS[41]	\checkmark	66.9	71.4	75.9	56.5	63.7	13.55
STM[45]	\checkmark	79.4	79.7	84.2	72.8	80.9	5.1

Table 8: Video segmentation results on Youtube-VOS dataset. Higher values are better. According to the evaluation protocol of the benchmark, we report performance separated into “seen” and “unseen” classes (“Seen” with respect to training set). [†] indicates results based on our reimplementation. The first- and second-best results on the unseen category are highlighted in red and blue, respectively.

distribution samples to gauge the model’s generalizability to more challenging, unseen, real-world scenarios. As seen from the last two columns, we rank second amongst all algorithms in unseen objects. In these unseen classes, we are even 3.9% higher than the DAVIS 2018 and YouTube-VOS 2018 video segmentation challenge winner, PreMVOS[41], a complex algorithm trained with multiple large manually labeled datasets. For fair comparison, we train our model only on the YouTube-VOS training set. We also re-train two most relevant self-supervised methods in the same manner as baselines. Even learning from only a subset of all classes, our model generalizes well to unseen classes, with a generalization gap (*i.e.* the performance difference between seen and unseen objects) near zero (0.4). This gap is much smaller than any of the baselines (avg = 11.5), suggesting

a unique advantage to most other algorithms trained with labels.

By training on large amounts of unlabeled videos, we learn an effective tracking representation without the need for any human annotations. This means that the learned network is not limited to a specific set of object categories (*i.e.* those in the training set), but is more likely to be a “universal feature representation” for tracking. Indeed, the only supervised algorithm that is comparable to our method in generalizability is OSVOS (2.7 vs. 0.4). However, OSVOS uses the first image from the testing sequence to perform costly domain adaptation, *e.g.* one-shot fine-tuning. In contrast, our algorithm requires no fine-tuning, which further demonstrates its zero-shot generalization capability.

Note our model also has a smaller generalization gap compared to other self-supervised methods as well. This further attests to the robustness of its learned features, suggesting that our improved reconstruction objective is highly effective in capturing general features.

6. Conclusion

In summary, we present a memory-augmented self-supervised model that enables accurate and generalizable pixel-level tracking. The algorithm is trained without any semantic annotation, and surpasses previous self-supervised methods on existing benchmarks by a significant margin, narrowing the gap with supervised methods. On unseen object categories, our model actually outperforms all but one existing methods that are trained with heavy supervision. As computation power grows and more high quality videos become available, we believe that self-supervised learning algorithms can serve as a strong competitor to their supervised counterparts for their flexibility and generalizability.

7. Acknowledgements

The authors would like to thank Andrew Zisserman for helpful discussions, Olivia Wiles, Shangzhe Wu, Sophia Koepke and Tengda Han for proofreading. Financial support for this project is provided by EPSRC Seebibyte Grant EP/M013774/1. Erika Lu is funded by the Oxford-Google DeepMind Graduate Scholarship.

References

- [1] Pulkit Agrawal, Joao Carreira, and Jitendra Malik. Learning to see by moving. In *Proc. ICCV*, 2015.
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *Proc. ICLR*, 2015.
- [3] Linchao Bao, Baoyuan Wu, and Wei Liu. Cnn in mrf: Video object segmentation via inference in a cnn-based higher-order spatio-temporal mrf. In *Proc. CVPR*, 2018.
- [4] Bert De Brabandere, Xu Jia, Tinne Tuytelaars, and Luc Van Gool. Dynamic filter networks. In *NIPS*, 2016.
- [5] T. Berry Brazelton, Mary Louise Scholl, and John S. Robey. Visual responses in the newborn. *Pediatrics*, 1966.
- [6] Sergi Caelles, Kevis-Kokitsi Maninis, Jordi Pont-Tuset, Laura Leal-Taixé, Daniel Cremers, and Luc Van Gool. One-shot video object segmentation. In *Proc. CVPR*, 2017.
- [7] Wu Chao-Yuan, Feichtenhofer Christoph, Fan Haoqi, He Kaiming, Krähenbühl Philipp, and Girshick Ross. Long-Term Feature Banks for Detailed Video Understanding. In *Proc. CVPR*, 2019.
- [8] Jingchun Cheng, Yi-Hsuan Tsai, Wei-Chih Hung, Shengjin Wang, and Ming-Hsuan Yang. Fast and accurate online video object segmentation via tracking parts. In *Proc. CVPR*, 2018.
- [9] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. Imagenet: A large-scale hierarchical image database. In *Proc. CVPR*, 2009.
- [11] Emily Denton and Vighnesh Birodkar. Unsupervised learning of disentangled representations from video. In *NIPS*, 2017.
- [12] J. Devlin, M.W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. 2018.
- [13] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The PASCAL Visual Object Classes Challenge 2009 (VOC2009) Results. <http://www.pascal-network.org/challenges/VOC/voc2009/workshop/index.html>, 2009.
- [14] Basura Fernando, Hakan Bilen, Efstratios Gavves, and Stephen Gould. Self-supervised video representation learning with odd-one-out networks. In *Proc. CVPR*, 2017.
- [15] Chuang Gan, Boqing Gong, Kun Liu, Hao Su, and Leonidas J. Guibas. Geometry guided convolutional neural networks for self-supervised video representation learning. In *Proc. CVPR*, 2018.
- [16] Alex Graves, Greg Wayne, and Ivo Danihelka. Neural Turing machines. *arXiv preprint arXiv:1410.5401*, 2014.
- [17] Tengda Han, Weidi Xie, and Andrew Zisserman. Video representation learning by dense predictive coding. In *1st International Workshop on Large-scale Holistic Video Understanding, ICCV*, 2019.
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. CVPR*, 2016.
- [19] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [20] Yuan-Ting Hu, Jia-Bin Huang, and Alexander G. Schwing. Maskrcnn: Instance level video object segmentation. In *NIPS*, 2017.
- [21] Yuan-Ting Hu, Jia-Bin Huang, and Alexander G. Schwing. Videomatch: Matching based video object segmentation. In *Proc. ECCV*, 2018.
- [22] Phillip Isola, Daniel Zoran, Dilip Krishnan, and Edward H. Adelson. Learning visual groups from co-occurrences in space and time. In *Proc. ICLR*, 2015.
- [23] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In *NIPS*, 2015.

- [24] Tomas Jakab, Ankush Gupta, Hakan Bilen, and Andrea Vedaldi. Conditional image generation for learning the structure of visual objects. In *NIPS*, 2018.
- [25] Dinesh Jayaraman and Kristen Grauman. Learning image representations tied to ego-motion. In *Proc. ICCV*, 2015.
- [26] Dinesh Jayaraman and Kristen Grauman. Slow and steady feature analysis: higher order temporal coherence in video. In *Proc. CVPR*, 2016.
- [27] Longlong Jing, Xiaodong Yang, Jingen Liu, and Yingli Tian. Self-supervised spatiotemporal feature learning by video geometric transformations. *arXiv preprint arXiv:1811.11387*, 2018.
- [28] Joakim Johnander, Martin Danelljan, Emil Brissman, Fahad Shahbaz Khan, and Michael Felsberg. A generative appearance model for end-to-end video object segmentation. In *Proc. CVPR*, 2019.
- [29] Anna Khoreva, Rodrigo Benenson, Eddy Ilg, Thomas Brox, and Bernt Schiele. Lucid data dreaming for multiple object tracking. In *arXiv preprint arXiv: 1703.09554*, 2017.
- [30] Anna Khoreva, Federico Perazzi, Rodrigo Benenson, Bernt Schiele, and Alexander Sorkine-Hornung. Learning video object segmentation from static images. In *arXiv preprint arXiv:1612.02646*, 2016.
- [31] A. Khoreva, A. Rohrbach, and B. Schiele. Video object segmentation with language referring expressions. In *Proc. ACCV*, 2018.
- [32] Dahun Kim, Donghyeon Cho, and In So Kweon. Self-supervised video representation learning with space-time cubic puzzles. In *AAAI*, 2018.
- [33] Shu Kong and Charless Fowlkes. Multigrid predictive filter flow for unsupervised learning on videos. *arXiv 1904.01693*, 2019, 2019.
- [34] Janet P. Kremenitzer, Herbert G. Vaughan, Diane Kurtzberg, and Kathryn Dowling. Smooth-pursuit eye movements in the newborn infant. *Child Development*, 1979.
- [35] Matej Kristan, Jiri Matas, Aleš Leonardis, Tomas Vojir, Roman Pflugfelder, Gustavo Fernandez, Georg Nebehay, Fatih Porikli, and Luka Čehovin. A novel performance evaluation methodology for single-target trackers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2016.
- [36] Ankit Kumar, Ozan Irsoy, Peter Ondruska, Mohit Iyyer, James Bradbury, Ishaan Gulrajani, Victor Zhong, Romain Paulus, and Richard Socher. Ask me anything: Dynamic memory networks for natural language processing. In *International conference on machine learning*, 2016.
- [37] Zihang Lai and Weidi Xie. Self-supervised learning for video correspondence flow. In *Proc. BMVC.*, 2019.
- [38] Hsin-Ying Lee, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Unsupervised representation learning by sorting sequences. In *Proc. ICCV*, 2017.
- [39] Xiaoxiao Li and Chen Change Loy. Video object segmentation with joint re-identification and attention-aware mask propagation. In *Proc. ECCV*, 2018.
- [40] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollar. Microsoft coco: Common objects in context. In *Proc. ECCV*, 2014.
- [41] Jonathon Luiten, Paul Voigtlaender, and Bastian Leibe. Premvos: Proposal-generation, refinement and merging for video object segmentation. In *Proc. ACCV*, 2018.
- [42] Kevis-Kokitsi Maninis, Sergi Caelles, Yuhua Chen, Jordi Pont-Tuset, Laura Leal-Taixé, Daniel Cremers, and Luc Van Gool. Video object segmentation without temporal information. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2018.
- [43] Ishan Misra, C. Lawrence Zitnick, and Martial Hebert. Shuffle and learn: Unsupervised learning using temporal order verification. In *Proc. ECCV*, 2016.
- [44] Seoung Wug Oh, Joon-Young Lee, Kalyan Sunkavalli, and Seon Joo Kim. Fast video object segmentation by reference-guided mask propagation. In *Proc. CVPR*, 2018.
- [45] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *Proc. ICCV*, 2019.
- [46] Aaron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *Proc. ICML*, 2016.
- [47] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Proc. CVPR*, 2016.
- [48] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbelaz, Alex Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv:1704.00675*, 2017.
- [49] Erik Reinhard, Michael Adhikhmin, Bruce Gooch, and Peter Shirley. Color transfer between images. *IEEE Computer graphics and applications*, 21(5):34–41, 2001.
- [50] Abigail See, Peter J Liu, and Christopher D Manning. Get to the point: Summarization with pointer-generator networks. *arXiv preprint arXiv:1704.04368*, 2017.
- [51] Sainbayar Sukhbaatar, Jason Weston, Rob Fergus, et al. End-to-end memory networks. In *Advances in neural information processing systems*, pages 2440–2448, 2015.
- [52] Jack Valmadre, Luca Bertinetto, Joao F. Henriques, Ran Tao, Andrea Vedaldi, Arnold Smeulders, Philip Torr, and Efstratios Gavves. Long-term tracking in the wild: A benchmark. In *Proc. ECCV*, 2018.
- [53] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In *NIPS*, 2017.
- [54] Carles Ventura, Miriam Bellver, Andreu Girbau, Amaia Salvador, Ferran Marques, and Xavier Giro-i Nieto. Rvos: End-to-end recurrent network for video object segmentation. In *Proc. CVPR*, 2019.
- [55] Carles Ventura, Miriam Bellver, Andreu Girbau, Amaia Salvador, Ferran Marques, and Xavier Giro-i Nieto. Rvos: End-to-end recurrent network for video object segmentation. In *Proc. CVPR*, June 2019.
- [56] Paul Voigtlaender, Yuning Chai, Florian Schroff, Hartwig Adam, Bastian Leibe, and Liang-Chieh Chen. Feelvos: Fast end-to-end embedding learning for video object segmentation. In *Proc. CVPR*, 2019.
- [57] Paul Voigtlaender and Bastian Leibe. Online adaptation of convolutional neural networks for video object segmentation. In *Proc. BMVC.*, 2017.
- [58] C. von Hofsten. Eyehand coordination in the newborn. *Developmental Psychology*, 1982.
- [59] Carl Vondrick, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy. Tracking emerges by colozing videos. In *Proc. ECCV*, 2018.
- [60] Ning Wang, Yibing Song, Chao Ma, Wengang Zhou, Wei

- Liu, and Houqiang Li. Unsupervised deep tracking. In *Proc. CVPR*, 2019.
- [61] Qiang Wang, Li Zhang, Luca Bertinetto, Weiming Hu, and Philip H.S. Torr. Fast online object tracking and segmentation: A unifying approach. In *Proc. CVPR*, 2019.
- [62] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proc. CVPR*, 2018.
- [63] Xiaolong Wang and Abhinav Gupta. Unsupervised learning of visual representations using videos. In *Proc. ICCV*, 2015.
- [64] Xiaolong Wang, Allan Jabri, and Alexei A. Efros. Learning correspondence from the cycle-consistency of time. In *Proc. CVPR*, 2019.
- [65] Ziqin Wang, Jun Xu, Li Liu, Fan Zhu, and Ling Shao. Ranet: Ranking attention network for fast video object segmentation. In *Proc. ICCV*, 2019.
- [66] Donglai Wei, Joseph Lim, Andrew Zisserman, and William T. Freeman. Learning and using the arrow of time. In *Proc. CVPR*, 2018.
- [67] Olivia Wiles, A. Sophia Koepke, and Andrew Zisserman. Self-supervised learning of a facial attribute embedding from video. In *Proc. BMVC.*, 2018.
- [68] Olivia Wiles, A. Sophia Koepke, and Andrew Zisserman. X2face: A network for controlling face generation using images, audio, and pose codes. In *Proc. ECCV*, 2018.
- [69] Zhu Xizhou, Wang Yujie, Dai Jifeng, Yuan Lu, and Wei Yichen. Flow-guided feature aggregation for video object detection. In *Proc. ICCV*, 2017.
- [70] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *Proc. ICML*, 2015.
- [71] Ning Xu, Linjie Yang, Yuchen Fan, Dingcheng Yue, Yuchen Liang, Jianchao Yang, and Thomas Huang. Youtubevos: A large-scale video object segmentation benchmark. *arXiv:1809.03327*, 2018.
- [72] Li Xueting, Liu Sifei, De Mello Shalini, Wang Xiaolong, Kautz Jan, and Yang Ming-Hsuan. Joint-task self-supervised learning for temporal correspondence. In *NeurIPS*, 2019.
- [73] Linjie Yang, Yanran Wang, Xuehan Xiong, Jianchao Yang, and Aggelos Katsaggelos. Efficient video object segmentation via network modulation. In *Proc. CVPR*, 2018.
- [74] Li Yao, Atousa Torabi, Kyunghyun Cho, Nicolas Ballas, Christopher Pal, Hugo Larochelle, and Aaron Courville. Describing videos by exploiting temporal structure. In *Proc. ICCV*, 2015.
- [75] Richard Zhang, Phillip Isola, and Alexei Efros. Colorful image colorization. In *Proc. ECCV*, 2016.

Appendix A. Network architecture

In the same way as CorrFlow [37], we use a modified ResNet-18 [18] architecture. Details of the network are illustrated in Table 9.

Stage	Output	Configuration
0	$H \times W$	Input image
conv1	$H/2 \times W/2$	$7 \times 7, 64, \text{stride } 2$
conv2	$H/2 \times W/2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$
conv3	$H/4 \times W/4$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$
conv4	$H/4 \times W/4$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$
conv5	$H/4 \times W/4$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$

Table 9: Network architecture. Residual Blocks are shown in brackets (a residually connected sequence of operations). See [18] for details.

Appendix B. Optimal memory size

In Figure 8, we explicitly show the effectiveness of increasing the number of reference frames in the memory bank, and confirm that a 5-frame memory is optimal for our task.

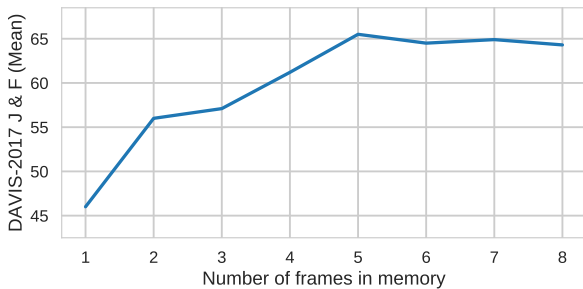


Figure 8: **Optimal memory size:** Here, we test a changing memory size of $n + m$: n short term memory and m long term memory, where n and m grow alternatively. The performance of our model initially increases as the number of frames in memory grows, eventually plateauing at 5 frames.

Appendix C. Analysis by attributes

We provide a more detailed accuracy list broken down by video attributes provided by the DAVIS benchmark [47] (listed in Table 10). The attributes illustrate the difficulties associated with each video sequence. Figure 9 contains the accuracies categorized by attribute. Several trends emerge: *first*, MAST outperforms all other self-supervised and unsupervised models by a large margin in all attributes. This shows that our model is robust to various challenges in dense tracking. *Second*, MAST obtains significant gains on occlusion-related video sequences (*e.g.* OCC, OV), suggesting that memory-augmentation is a key enabler for high-quality tracking: retrieving occluded objects from pre-

vious frames is very difficult without memory augmentation. *Third*, in videos involving background clutter, *i.e.* background and foreground share similar colors, MAST obtains a relatively small improvement over previous methods. We conjecture this bottleneck could be caused by a shared photometric loss; thus a different loss type (*e.g.* based on texture consistency) could further improve the result.

ID	Description	ID	Description
AC	Appearance Change	IO	Interacting Objects
BC	Background Clutter	LR	Low Resolution
CS	Camera-Shake	MB	Motion Blur
DB	Dynamic Background	OCC	Occlusion
DEF	Deformation	OV	Out-of-view
EA	Edge Ambiguity	ROT	Rotation
FM	Fast-Motion	SC	Shape Complexity
HO	Heterogeneous Object	SV	Scale-Variation

Table 10: List of video attributes provided in the DAVIS benchmark. We break down the validation accuracy according to the attribute list.

Appendix D. YouTube-VOS 2019 dataset

We also evaluate MAST and two other self-supervised methods on YouTube-VOS 2019 validation dataset. The numerical results are reported in Table 11. Augmenting on the 2018 version, the 2019 version contains more videos and object instances. We observe similar trend as reported in the main paper (*i.e.* significant improvement and lower generalization gap).

Method	Sup.	Overall \uparrow	Seen		Unseen		Gen. Gap \downarrow
			$\mathcal{J} \uparrow$	$\mathcal{F} \uparrow$	$\mathcal{J} \uparrow$	$\mathcal{F} \uparrow$	
Vid. Color.[59] [†]	\times	39.0	43.3	38.2	36.6	37.5	3.7
CorrFlow[37]	\times	47.0	51.2	46.6	44.5	45.9	3.7
MAST (Ours)	\times	64.9	64.3	65.3	61.5	68.4	0.15

Table 11: Video segmentation results on Youtube-VOS 2019 dataset. Higher values are better. [†] indicates results based on our reimplementation.

Appendix E. More qualitative results

As shown in Figure 10, we provide more qualitative results exhibiting some of difficulties in the tracking task. These difficulties include tracking multiple similar objects (multi-instance tracking often fails by conflating similar objects), large camera shake (objects may have motion blur), inferring unseen object pose of objects, and so on. As shown in the figure, MAST handles these difficulties well.

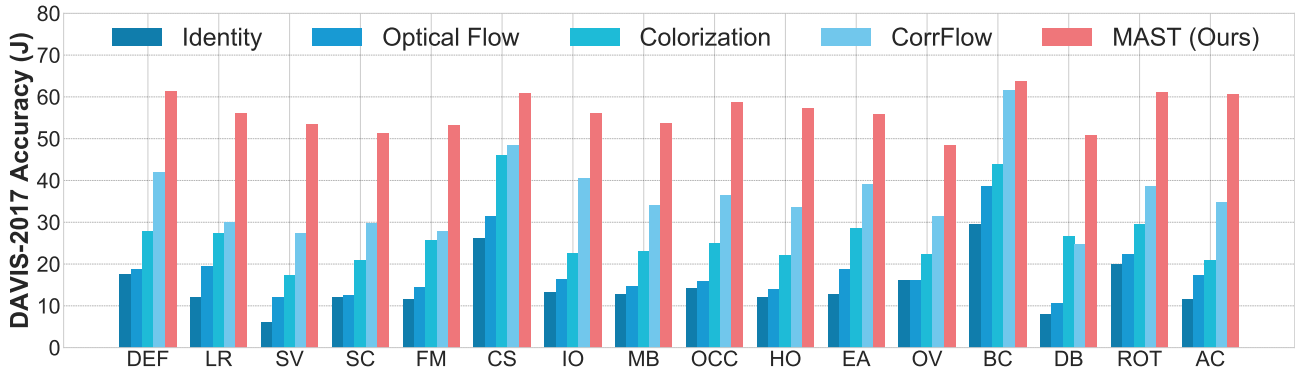


Figure 9: **Accuracy broken down by attribute:** MAST outperforms previous self-supervised methods by a significant margin on all attributes, demonstrating the robustness of our model.

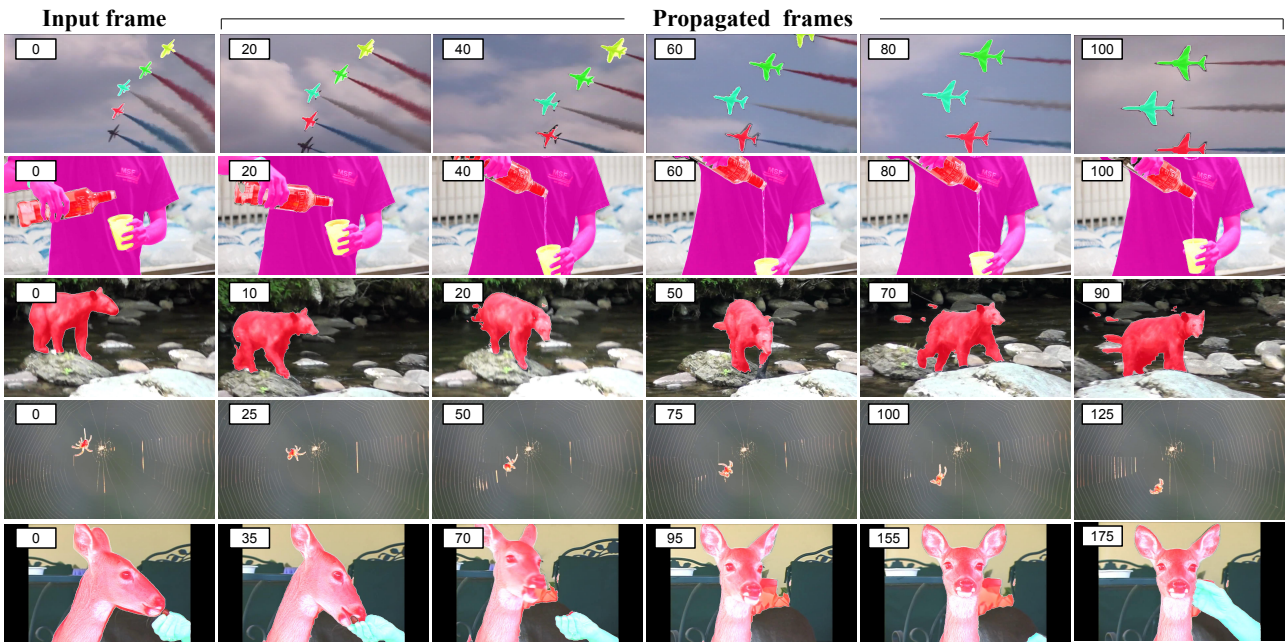


Figure 10: **More qualitative results** from our *self-supervised dense tracking model* on the YouTube-VOS dataset. The number on the top left refers to the frame number in the video. Row 1: Tracking multiple similar objects with scale change. Row 2: Occlusions and out-of-scene objects (hand, bottle, and cup). Row 3: Large camera shake. Row 4: Small object with fine details. Row 5: Inferring unseen pose of the deer; out-of-scene object (hand).


Statement of Authorship for joint/multi-authored papers for PGR thesis

To appear at the end of each thesis chapter submitted as an article/paper

The statement shall describe the candidate's and co-authors' independent research contributions in the thesis publications. For each publication there should exist a complete statement that is to be filled out and signed by the candidate and supervisor (**only required where there isn't already a statement of contribution within the paper itself**).

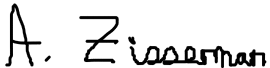
Title of Paper	MAST: A Memory-Augmented Self-Supervised Tracker
Publication Status	<input type="checkbox"/> Published <input type="checkbox"/> Accepted for Publication <input type="checkbox"/> Submitted for Publication <input type="checkbox"/> Unpublished and unsubmitted work written in a manuscript style
Publication Details	Lai, Z., Lu, E., & Xie, W. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020.

Student Confirmation

Student Name:	Erika Lu		
Contribution to the Paper	<ul style="list-style-type: none">• Development of idea• Writing the paper		
Signature 	Date	25 th August, 2021	

Supervisor Confirmation

By signing the Statement of Authorship, you are certifying that the candidate made a substantial contribution to the publication, and that the description described above is accurate.

Supervisor name and title:	Professor Andrew Zisserman		
Supervisor comments			
Signature 	Date	9/9/2021	

This completed form should be included in the thesis, at the end of the relevant chapter.

6

Self-Supervised Video Object Segmentation by Motion Grouping

This work will be presented at the IEEE International Conference on Computer Vision (ICCV), 2021.

Self-supervised Video Object Segmentation by Motion Grouping

Charig Yang Hala Lamdouar Erika Lu
 Andrew Zisserman Weidi Xie

Visual Geometry Group, University of Oxford

{charig, lamdouar, erika, az, weidi}@robots.ox.ac.uk

<https://charigyang.github.io/motiongroup/>

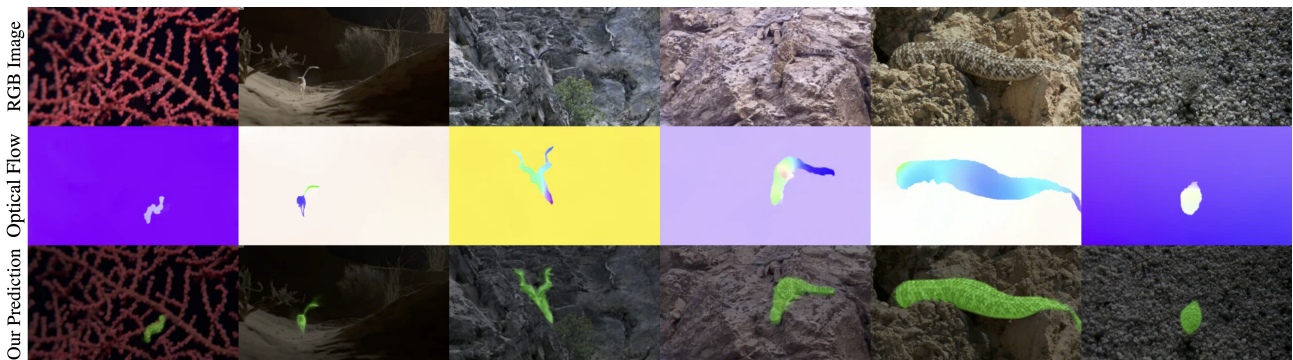


Figure 1. **Segmenting camouflaged animals.** Motion plays a critical role in augmenting the capability of our visual system for perceptual grouping in complex scenes – for example, in these sequences (MoCA dataset [41]), the visual appearance (RGB images) is clearly uninformative. We propose a self-supervised approach to segment objects using *only* motion, *i.e.* optical flow. From top to bottom rows, we show the video frames, optical flow between consecutive frames, and the segmentation produced by our approach.

Abstract

Animals have evolved highly functional visual systems to understand motion, assisting perception even under complex environments. In this paper, we work towards developing a computer vision system able to segment objects by exploiting motion cues, *i.e.* motion segmentation. We make the following contributions: First, we introduce a simple variant of the Transformer to segment optical flow frames into primary objects and the background. Second, we train the architecture in a self-supervised manner, *i.e.* without using any manual annotations. Third, we analyze several critical components of our method and conduct thorough ablation studies to validate their necessity. Fourth, we evaluate the proposed architecture on public benchmarks (DAVIS2016, SegTrackv2, and FBMS59). Despite using only optical flow as input, our approach achieves superior or comparable results to previous state-of-the-art self-supervised methods, while being an order of magnitude faster. We additionally evaluate on a challenging camouflage dataset (MoCA), significantly outperforming the other self-supervised approaches, and comparing favourably to the top supervised

approach, highlighting the importance of motion cues, and the potential bias towards visual appearance in existing video segmentation models.

1. Introduction

When looking around the world, we effortlessly perceive a complex scene as a set of distinct objects. This phenomenon is referred to as *perceptual grouping* – the process of organizing the incoming visual information – and is usually considered a fundamental cognitive ability that enables understanding and interacting with the world efficiently. How do we accomplish such a remarkable perceptual achievement, given that the visual input is, in a sense, just a spatial distribution of variously colored individual points/pixels? In 1923, Wertheimer [80] first introduced the *Gestalt principles* with the goal of formulating the underlying causes by which sensory data is organized into groups, or Gestalten. The principles are much like heuristics with “a bag of tricks” [59] that the visual system may exploit for grouping, for example, proximity, similarity, closure, continuation, common fate, *etc.*

In computer vision, *perceptual grouping* is often closely related to the problem of segmentation, *i.e.* extracting the objects with arbitrary shape (pixel-wise labels) from cluttered scenes. In the recent literature of semantic or instance segmentation, tremendous progress has been made by training deep neural networks on image or video datasets. While it is exciting to see machines with the ability to detect, segment, and classify objects in images or video frames, training such segmentation models through supervised learning requires massive human annotation, and consequently limits their scalability. Even more importantly, the assumption that objects can be well-identified by their appearance alone in static frames is often an oversimplification – objects are not always visually distinguishable from their background environment. For instance when trying to discover camouflaged animals/objects from the background (Figure 1), extra cues, such as motion or sound, are usually required.

Among the numerous cues, motion is usually simple to obtain as it can be implicitly generated from unlabeled videos. In this paper, we aim to exploit such cues for object segmentation in a self-supervised manner, *i.e.* zero human annotation is required for training. At a high level, we aim to exploit the common fate principle, with the basic assumption being that **elements tend to be perceived as a group if they move in the same direction at the same rate (have similar optical flow)**. Specifically, we tackle the problem by training a generative model that decomposes the optical flow into foreground (object) and background layers, describing each as a homogeneous field, with discontinuities occurring only between layers. We adopt a variant of the Transformer [3], with the self-attention being replaced by slot attention [46], where iterative grouping and binding have been built into the architecture. With some critical architectural changes, we show that pixels undergoing similar motion are grouped together and assigned to the same layer.

To summarize, we make the following contributions: *first*, we introduce a simple architecture for video object segmentation by exploiting motions, using only optical flow as input. *Second*, we propose a self-supervised proxy task that is used to train the architecture without any manual supervision. *Third*, we conduct thorough ablation studies on the components that are key to the success of our architecture, such as a consistency loss on optical flow computed from various frame gaps. *Fourth*, we evaluate the proposed architecture on public benchmarks (DAVIS2016 [55], SegTrackv2 [42], and FBMS59 [56]), outperforming previous state-of-the-art self-supervised models, with comparable performance to the supervised approaches. Moreover, we also evaluate on a camouflage dataset (MoCA [41]), demonstrating a significant performance improvement over the other self- and supervised approaches, highlighting the importance of motion cues, and the potential bias towards visual appearance in existing video segmentation models.

2. Related Work

Video object segmentation has been a longstanding task in computer vision, which involves segmenting the pixels (or edges) belonging to an image into groups (e.g. objects). In recent literature [5, 10, 12, 16, 24, 25, 31, 34, 35, 39, 40, 51, 52, 53, 54, 57, 57, 71, 73, 74, 75, 77, 84, 84, 87], two protocols have attracted increasing interest from the vision community, namely, semi-supervised video object segmentation (**semi-supervised VOS**), and unsupervised video object segmentation (**unsupervised VOS**). The former aims to re-localize one or multiple targets that are specified in the first frame of a video with pixel-wise masks, and the latter considers automatically separating the object of interest (usually the most salient one) from the background in a video sequence. Despite being called **unsupervised VOS**, in practice, the popular methods to address such problems extensively rely on supervised training, for example, by using two-stream networks [16, 31, 54, 71] trained on large-scale external datasets. As an alternative, in this work, we consider a *completely unsupervised* problem, where no manual annotation is used for training whatsoever.

Motion segmentation shares some similarity with unsupervised VOS, but focuses on discovering *moving* objects. In the literature, [10, 36, 52, 63, 83] consider clustering the pixels with similar motion patterns; [16, 70, 71] train deep networks to map the motions to segmentation masks. Another line of work has tackled the problem by explicitly leveraging the independence, in the flow field, between the moving object and its background. For instance, [86] proposes an adversarial setting, where a generator is trained to produce masks, altering the input flow, such that the inpainter fails to estimate the missing information. In [6, 7, 41], the authors propose to highlight the independently moving object by compensating for the background motion, either by registering consecutive frames, or explicitly estimating camera motion. In constrained scenarios, such as the autonomous driving domain, [61] proposes to jointly optimize depth, camera motion, optical flow and motion segmentation.

Optical flow computation is one of the fundamental tasks in computer vision. Deep learning methods allow efficient computation of optical flow, both in supervised learning on synthetic data [68, 69], and in the self-supervised [44, 45] setting. Flow has been useful for a wide range of problems, occasionally even used in lieu of appearance cues (RGB images) for tracking [62], pose estimation [19], representation learning [50], and motion segmentation [10].

Transformer architectures have proven extremely adept at modeling long-term relationships within an input sequence via attention mechanisms. Originally used for language tasks [3, 9, 18], they have since been adapted for solving popular computer vision problems such as

image classification [20], generation [14, 60], video understanding [4, 27, 79], object detection [13], and zero-shot classification [58]. In this work, we take inspiration from a specific variant of self-attention, namely slot attention [46], which was demonstrated to be effective for learning object-centric representations on synthetic datasets, *e.g.* CLEVR.

Layered representations were originally proposed by Wang and Adelson [76] to represent a video as a composition of layers with simpler motions. Since then, layered representations of images and videos have been widely adopted in computer vision [8, 33, 38, 85, 90], often to estimate optical flow [66, 67, 81, 82]. More recently, deep learning-based layer decomposition methods have been used to infer depth for novel view synthesis [65, 88], separate reflections and other semi-transparent effects [1, 2, 26, 47], or perform foreground/background estimation [26]. These works operate on RGB inputs and produce RGB layers, whereas we propose a layered decomposition of optical flow inputs for unsupervised moving object discovery.

Object-centric representations interpret scenes with “objects” as the basic building block (instead of individual pixels), which is considered an essential step towards human-level generalization. There is a rich literature on this topic, for example, in [29], the IODINE model uses iterative variational inference to infer a set of latent variables recurrently, with each representing one object in an image. Similarly, MONet [11] and GENESIS [21] also adopt the multiple encode-decode steps. In contrast, [46] proposes Slot Attention, which enables single step encoding-decoding with iterative attention. However, all of the works mentioned above have only shown applications for synthetic datasets, *e.g.* CLEVR [32]. In this paper, we are the first to demonstrate its use for object segmentation of realistic videos by exploiting motion, where the challenging nuances in visual appearance (*e.g.* the complex background textures) have been removed.

3. Method

Our goal is to take an input optical flow frame and predict a segment containing the moving object. We aim to train this model in a self-supervised manner, specifically, we adopt an autoencoder-like framework. Our model outputs two layers: one representing the background, and the other for one or more moving objects in the foreground, as well as their opacity layers (weighted masks). Formally, we have:

$$\{\hat{I}_{t \rightarrow t+n}^i, \alpha_{t \rightarrow t+n}^i\}_{i=1}^N = \Phi(I_{t \rightarrow t+n}) \quad (1)$$

where $I_{t \rightarrow t+n}$ refers to the t to $t+n$ input flow (backward flow when $n < 0$), $\Phi(\cdot)$ is the parametrized model, $\hat{I}_{t \rightarrow t+n}^i$ is the i th layer reconstruction, $\alpha_{t \rightarrow t+n}^i$ is its mask,

and $N = 2$ is the number of layers (foreground and background). These layers can then be composited linearly to reconstruct the input image $I_{t \rightarrow t+n}$:

$$\hat{I}_{t \rightarrow t+n} = \sum_{i=1}^N \alpha_{t \rightarrow t+n}^i \hat{I}_{t \rightarrow t+n}^i \quad (2)$$

3.1. Flow Segmentation Architecture

For simplicity, we first consider the case of a single flow field as input (depicted in the top part of Figure 2). The entire model consists of three components: (1) a CNN encoder to extract a compact feature representation, (2) an iterative binding module with learnable queries that plays a similar role as soft clustering, *i.e.* assigning each pixel to one of motion groups, and (3) a CNN decoder that individually decodes each query to full resolution layer outputs (where thresholding the alpha channel yields the predicted segment).

CNN Encoder. We first pass the precomputed optical flow field between two frames, $I_{t \rightarrow t+n} \in \mathcal{R}^{3 \times H_0 \times W_0}$, to a CNN encoder Φ_{enc} , which outputs a lower-resolution feature map:

$$F_{t \rightarrow t+n} = \Phi_{\text{enc}}(I_{t \rightarrow t+n}) \in \mathcal{R}^{D \times H \times W} \quad (3)$$

where H_0, W_0 and H, W refer to the spatial dimensions of the input and output feature maps respectively. Note that, we convert the optical flow field into a three-channel image according to the traditional method specified in the optical flow literature [68].

Iterative Binding. The iterative binding module Φ_{bind} aims to group image regions into single entities based on their similarities in motion, *i.e.* pixels moving in the same direction at the same rate should be grouped together. Intuitively, such a binding process requires a data-dependent parameter updating mechanism, iteratively adapting and enriching the model, gradually including more pixels undergoing similar motions.

To accomplish this task, we adopt a simple variant of slot attention [46], where instead of Gaussian-initialized slots, we use learnable query vectors. Slot attention has recently shown remarkable performance for object-centric representation learning, where the query vectors compete to explain parts of the inputs via a softmax-based attention mechanism, and the representations of these slots are iteratively updated with a recurrent update function. In our case of motion segmentation, ideally, the final representation in each query vector separately encodes the moving object or the background, which can then be decoded and combined to reconstruct the input flow fields.

Formally, our inputs to Φ_{bind} are feature maps $F_{t \rightarrow t+n}$ and two learnable queries (representing foreground and

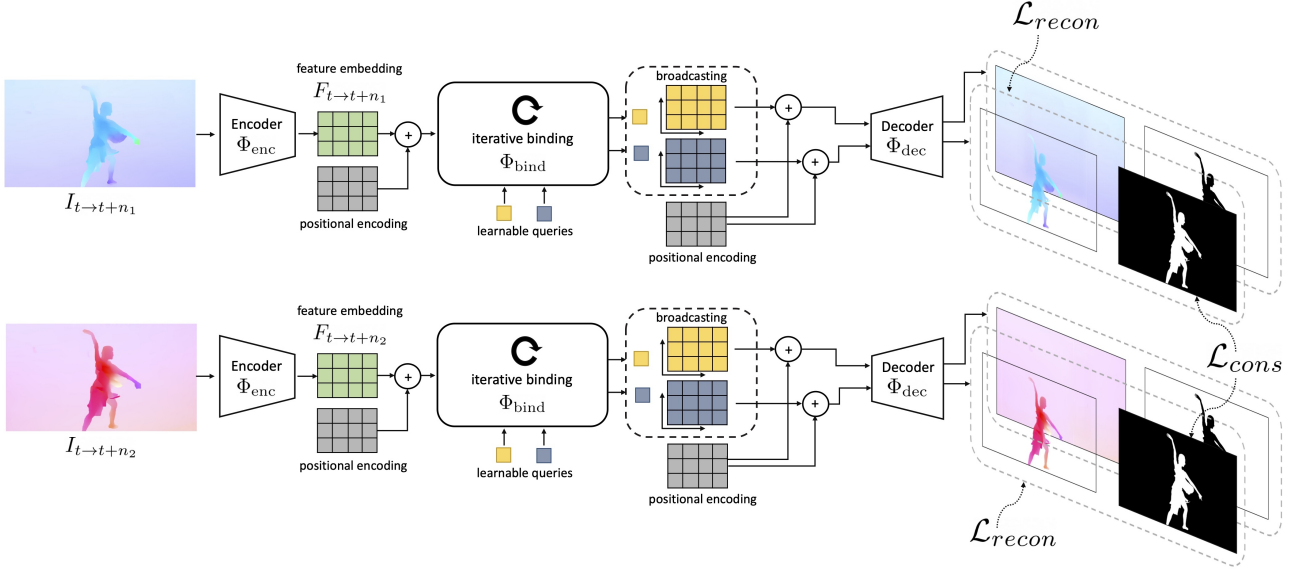


Figure 2. **Pipeline.** Our model takes optical flow $I_{t \rightarrow t+n}$ as input, and outputs a set of reconstruction and opacity layers. Specifically, it consists of three components: feature encoding, iterative binding, and decoding to layers, which are combined to reconstruct the input flow. To resolve motion ambiguities (small motion), or noise in optical flow, consistency between two flow fields computed under different frame gaps is enforced during training. At inference time, only the top half of the figure is used to predict masks from a single-step flow.

background) $Q \in \mathcal{R}^{D \times 2}$. Learnable spatial positional encodings are summed with $F_{t \rightarrow t+n}$; with some abuse of notation, we still refer to this sum as $F_{t \rightarrow t+n}$. We use *three* different linear transformations to generate the *query*, *key* and *value*: $q \in \mathcal{R}^{D \times 2}$, $k, v \in \mathcal{R}^{D \times HW}$,

$$q, k, v = W^Q \cdot Q, W^K \cdot F_{t \rightarrow t+n}, W^V \cdot F_{t \rightarrow t+n} \quad (4)$$

where $W^Q, W^K, W^V \in \mathcal{R}^{D \times D}$.

In contrast to the standard Transformer [3], the coefficients in slot attention are normalized over all slots. This choice of normalization introduces competition between the slots to explain parts of the input, and ensures each pixel is assigned to a query vector:

$$\text{attn}_{i,j} := \frac{e^{M_{i,j}}}{\sum_l e^{M_{i,l}}} \quad (5)$$

$$M := \frac{1}{\sqrt{D}} k^T \cdot q, \quad \text{attn} \in \mathcal{R}^{HW \times 2}$$

To aggregate the input values to their assigned query slot, a weighted mean is used as follows:

$$U := v \cdot A \in \mathcal{R}^{D \times 2} \quad (6)$$

where, $A_{i,j} := \frac{\text{attn}_{i,j}}{\sum_l \text{attn}_{i,l}}$

To maintain a smooth update of the query slots Q , the aggregated vectors U are fed into a recurrent function, parametrized with Gated Recurrent Units (GRU),

$$Q := \text{GRU}(\text{inputs} = U, \text{states} = Q) \quad (7)$$

This whole binding process is then iterated T times. The pseudocode can be found in the Supplementary Material.

CNN Decoder. The CNN decoder Φ_{dec} individually decodes each of the slots to outputs of original resolution ($\{\hat{I}_{t \rightarrow t+n}^i, \alpha_{t \rightarrow t+n}^i\} \in \mathcal{R}^{4 \times H_0 \times W_0}$), which includes an (unnormalized) single-channel alpha mask and the reconstructed flow fields. Specifically, the input to the decoder is the slot vector broadcasted onto a 2D grid augmented with a learnable spatial positional encoding.

Reconstruction. Once each slot has been decoded, we apply softmax to the alpha masks across the slot dimension, and use them as mixture weights to obtain the reconstruction $\hat{I}_{t \rightarrow t+n}$ (Eq. 2). Our reconstruction loss is an L2 loss between the input and reconstructed flow,

$$\mathcal{L}_{\text{recon}} = \frac{1}{\Omega} \sum_{p \in \Omega} |I_{t \rightarrow t+n}(p) - \hat{I}_{t \rightarrow t+n}(p)|^2 \quad (8)$$

where p is the pixel index, and Ω is the entire spatial grid.

Entropy Regularization. We impose a pixel-wise entropy regularisation on inferred masks:

$$\mathcal{L}_{\text{entr}} = \frac{1}{\Omega} \sum_{p \in \Omega} (-\alpha_{t \rightarrow t+n}^1(p) \log \alpha_{t \rightarrow t+n}^1(p) - \alpha_{t \rightarrow t+n}^0(p) \log \alpha_{t \rightarrow t+n}^0(p)) \quad (9)$$

This loss is zero when the alpha channels are one-hot, and maximum when they are of equal probability. Intuitively, this helps encourage the masks to be binary, which aligns

with our goal in obtaining segmentation masks.

Instance Normalisation. In the case of motion segmentation, objects can only be detected if they undergo an independent motion from the camera; thus previous work attempts to compensate for camera motion [6, 41]. We are inspired by these ideas, but instead of explicitly estimating homography or camera motion, we take a poor-man’s approach by simply using Instance Normalisation (IN) [72] in the CNN encoder and decoder, which normalizes each channel in each training sample independently. Intuitively, the *mean* activation tends to be dominated by the motions in the large homogeneous region, which is usually the background. This normalization, in combination with ReLU activations, helps gradually separate the background motion from the foreground motions.

3.2. Self-supervised Temporal Consistency Loss

The segmentation computed for the current frame should be identical irrespective of whether the ‘second’ frame is consecutive, or earlier or later in time. We harness this constraint to form a self-supervised temporal consistency loss by first defining a set of ‘second’ frames, and then requiring consistency between their pairwise predictions. We describe the set first, followed by the loss.

Multi-step flow. As objects may be static for some frames, we make our predictions more robust by leveraging observations from multiple timesteps. We consider the flow fields computed from various temporal gaps as an input set, *i.e.* $\{I_{t \rightarrow t+n_1}, I_{t \rightarrow t+n_2}\}$, $n_1, n_2 \in \{-2, -1, 1, 2\}$, and use a permutation invariant consistency loss to encourage the model to predict the same foreground/background segmentation for all flow fields in the set.

Consistency loss. We randomly sample two flow fields from the input set and pass them through the model ($\Phi(\cdot)$), outputting the flow reconstruction and alpha masks for each. As the reconstruction loss is commutative, it is not guaranteed that the same slot will always output the background layer; therefore, we use a permutation-invariant consistency loss, *i.e.* only backpropagating through the lowest-error permutation:

$$\mathcal{L}_{cons} = \frac{1}{\Omega} \min \left(\sum_{p \in \Omega} |\alpha_{t \rightarrow t+n_1}^1(p) - \alpha_{t \rightarrow t+n_2}^1(p)|^2, \sum_{p \in \Omega} |\alpha_{t \rightarrow t+n_1}^1(p) - \alpha_{t \rightarrow t+n_2}^0(p)|^2 \right)$$

Note that, this consistency enforcement only occurs during training. At inference time, a single-step flow is used, as shown in the top half of Figure 2.

Total Loss. The total loss for training the architecture is:

$$\mathcal{L}_{total} = \gamma_r \mathcal{L}_{recon} + \gamma_c \mathcal{L}_{cons} + \gamma_e \mathcal{L}_{entr} \quad (10)$$

we use $\gamma_r = 10^2$, $\gamma_c = 10^{-2}$ and $\gamma_e = 10^{-2}$. However, we found the model is fairly robust to these hyper-parameters.

3.3. Discussion

Differences from slot attention. Slot attention was originally introduced for self-supervised object segmentation for RGB images [46], and its usefulness was demonstrated on synthetic data (CLEVR [32]), where objects are primitive shapes with simple textures. However, this assumption is unlikely to hold in the case of natural images or videos, making it challenging to generalize such object-centric representations.

In this work, we build on the insight that although objects in images may not be naturally textureless, their motions typically are. Hence, we develop the self-supervised object segmentation model by exploiting their optical flows, where the nuance in visual appearance is discarded, thus not restricted to simple synthetic cases. As an initial trial, we experimented with the same setting as [46], where query vectors are sampled from a Gaussian distribution; however, we were unable to train it. Instead, we use a learnable embedding here, which we highlight as one of the architectural changes critical to our model’s success. Other critical changes include instance normalization and temporal consistency, which we demonstrate in ablations in Section 5.1.

Why does it work for motion segmentation? The proposed idea can be seen as training a generative model to segment the flow fields. With the layered formulation, reconstruction is limited to be a simple *linear* composition of layer-wise flow, decoded from a single slot vector.

Conceptually, this design has effectively introduced a representational bottleneck, encouraging each slot vector to represent minimal information, *i.e.* homogeneous motion, and with minimal redundancy (mutual information) between slots. All these properties make such an architecture well-suited to the task of segmenting objects undergoing independent motions.

4. Experimental Setup

In this section, we describe the datasets and evaluation metrics, followed by implementation details.

4.1. Datasets

We benchmark on four different datasets that are commonly used for unsupervised video object segmentation.

DAVIS2016 [55] contains a total of 50 sequences (30 for training and 20 for validation), depicting diverse moving objects such as animals, people, and cars. The dataset contains 3455 1080p frames with pixel-wise annotations at 480p for the predominantly moving object.

SegTrackv2 [42] contains 14 sequences and 976 annotated frames. Each sequence contains 1-6 moving objects, and presents challenges including motion blur, appearance change, complex deformation, occlusion, slow motion, and interacting objects.

FBMS59 [56] consists of 59 sequences and 720 annotated frames (every 20th frame is annotated), which vary greatly in image resolution. Sequences involve multiple moving objects, some of which may be static for periods of time.

Moving Camouflaged Animals (MoCA) [41] contains 141 HD video sequences, depicting 67 kinds of camouflaged animals moving in natural scenes. Both temporal and spatial annotations are provided in the form of tight bounding boxes for every 5th frame. Using the provided motion labels (locomotion, deformation, static), we filter out videos with predominantly no locomotion, resulting in 88 video sequences and 4803 frames (this subset will be released).

4.2. Evaluation Metrics

We use two different evaluation metrics, depending on whether pixel-wise annotations or bounding boxes are provided.

Segmentation (Jaccard). For DAVIS2016, SegTrackv2 and FBMS59, pixelwise segmentation is provided; thus we report the standard metric, region similarity (\mathcal{J}), computing the mean over the test set. For FBMS59 and SegTrackv2, we follow the common practice [30, 86] and combine multiple objects as one single foreground.

Localization (Jaccard & Success Rate). As the MoCA dataset provides only bounding box annotations, we evaluate for the detection task and report results in the form of detection success rate [22, 43], for varying IoU thresholds ($\tau \in \{0.5, 0.6, 0.7, 0.8, 0.9\}$).

4.3. Implementation Details

We evaluate three different approaches for computing optical flow, namely, PWC-Net [68], RAFT [69] and ARFlow [44]; the first two are supervised, while the latter is self-supervised. We extract the optical flow at the original resolution of the image pairs, with the frame gaps $n \in \{-2, -1, 1, 2\}$ for all datasets, except for FBMS59, where we use $n \in \{-6, -3, 3, 6\}$ to compensate for small motion. To generate inputs to the network for training, the flows are resized to 128×224 (and scaled accordingly), converted to 3-channel images with the standard visualization used for optical flow, and normalized to $[-1, 1]$.

In the iterative binding module ($\Phi_{\text{bind}}(\cdot)$), we use two learnable query vectors (as we consider the case of segmenting a single moving object from the background), and choose $T = 5$ iterations (as explained in Section 3.1). We adopt a simple VGG-style network for the CNN encoder and decoder, with 3 blocks each; in each block, we use 2

sets of convolutions, instance normalization and ReLU activation. We train with a batch size of 64 images and use the Adam optimizer [37] with an initial learning rate of 5×10^{-4} , decreasing every 80k iterations. The exact architecture description and training schedule can be found in the Supplementary Material, and the submitted source code.

5. Results

In this section, we compare primarily with a top-performing approach trained without manual annotations – Contextual Information Separation (CIS [86]). However, as the architecture, input resolution, modality and post-processing are all different, we try our best to conduct the comparison as fairly as possible. Note that benchmarks are evaluated at full resolution by simply upsampling the predicted masks.

In order to guarantee a fair comparison of the model, we disregard the gains from post-processing, which typically include averaging an ensemble of multiple crops and flow steps, temporal smoothing, and CRFs, at the expense of runtime. Here, we consider only raw framewise predictions with single-step flow at over 80fps, which is more likely to be the application of motion segmentation in practice.

5.1. Ablation Studies

We conduct all ablation studies on DAVIS2016, and vary one variable each time, as shown in Table 1.

Choice of Optical Flow Algorithm. With the same flow extraction method (PWC-Net), our proposed model (Ours-A) outperforms CIS by about 4.5 points on mean Jaccard (\mathcal{J}), and using improved optical flow (RAFT) provides further performance gains. We therefore use RAFT from hereon.

Instance Normalization and Grouping. We observe two phenomena: *first*, when holding constant the number of grouping iterations T (3 or 5), models trained with instance normalization perform consistently better; *second*, iterative grouping with $T = 5$ is better than that trained with $T = 3$. However, at $T = 8$, the model did not converge in the same number of training steps, and thus we do not include it in the table. For the remainder of the experiments, we use instance normalization and $T = 5$ to balance performance and training time.

Consistency and Entropy Regularization. While comparing Ours-B and Ours-I, we observe that the performance degrades significantly without the temporal consistency loss, and that the entropy regularization is also important, as shown by Ours-B and Ours-H.

5.2. Comparison with State-of-the-art

We show comparison results for different datasets in Table 2. On DAVIS2016, we improve upon the state-of-the-art

Model	Flow	IN	T	\mathcal{L}_e	\mathcal{L}_c	DAVIS ($\mathcal{J} \uparrow$)
CIS [86]	PWC-Net	–	–	–	–	59.2
Ours-A	PWC-Net	✓	5	✓	✓	63.7
Ours-B	RAFT	✓	5	✓	✓	68.3
Ours-C	ARFlow	✓	5	✓	✓	53.2
Ours-D	RAFT	✓	3	✓	✓	65.8
Ours-E	RAFT	✗	3	✓	✓	63.3
Ours-F	RAFT	✗	5	✓	✓	64.5
Ours-G	RAFT	✓	5	✗	✗	48.0
Ours-H	RAFT	✓	5	✗	✓	60.3
Ours-I	RAFT	✓	5	✓	✗	51.2

Table 1. **Ablation studies** on flow extraction methods, instance normalization (IN), grouping iterations (T), entropy regularization (\mathcal{L}_e) and set consistency (\mathcal{L}_c).

Model	Sup.	RGB	Flow	Res.	DAVIS16 ($\mathcal{J} \uparrow$)	STv2 ($\mathcal{J} \uparrow$)	FBMS59 ($\mathcal{J} \uparrow$)	Runtime (sec \downarrow)
SAGE [78]	✗	✓	✓	–	42.6	57.6	61.2	0.9s
NLC [23]	✗	✓	✓	–	55.1	67.2	51.5	11s
CUT [36]	✗	✓	✓	–	55.2	54.3	57.2	103s
FTS [54]	✗	✓	✓	–	55.8	47.8	47.7	0.5s
CIS [86]	✗	✓	✓	192 × 384	59.2 (71.5)	45.6 (62.0)	36.8 (63.5)	0.1s (11s)
Ours	✗	✗	✓	128 × 224	68.3	58.6	53.1	0.012s
SFL [15]	✓	✓	✓	854 × 480	67.4	–	–	7.9s
FSEG [30]	✓	✓	✓	854 × 480	70.7	61.4	68.4	–
LVO [71]	✓	✓	✓	–	75.9	57.3	65.1	–
ARP [64]	✓	✓	✓	–	76.2	57.2	59.8	74.5s
COSNet [48]	✓	✓	✗	473 × 473	80.5	–	75.6	–
MATNet [89]	✓	✓	✓	473 × 473	82.4	–	–	0.55s
3DC-Seg [49]	✓	✓	✓	854 × 480	84.3	–	–	0.84s

Table 2. **Full comparison on moving object segmentation** (unsupervised video segmentation). We consider three popular datasets, DAVIS2016, SegTrack-v2 (STv2), and FBMS59. Models above the horizontal dividing line are trained without using *any* manual annotation, while models below are pre-trained on image or video segmentation datasets, *e.g.* DAVIS, YouTube-VOS, thus requiring ground truth annotations at training time. Numbers in parentheses denote the additional usage of significant post-processing, *e.g.* multi-step flow, multi-crop, temporal smoothing, CRFs.

for unsupervised methods (CIS) by a large margin (+9.1%). As shown in Figure 5, although we do not use any pixel-level annotations during training, our method is nearing the performance of supervised models trained on thousands of images. In addition, we argue that, motion segmentation in realistic scenarios, *e.g.* by predator or prey, is likely to require fast processing. Our model operates on small resolution (potentially sacrificing some accuracy) with over 80fps.

For SegTrackv2 and FBMS59, they occasionally include multiple objects in a single video, and only a subset of them are moving, making it challenging to spot all objects using flow-only input, but we achieve competitive performance nonetheless. We discuss this limitation below.

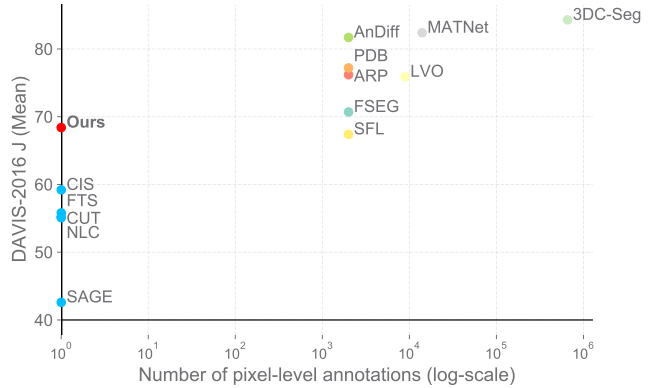


Figure 3. **Comparison on DAVIS2016**. Note that, supervised approaches may use models pretrained on ImageNet [17], but here we only count number of images with pixel-wise annotations.

5.3. Camouflage Breaking

In addition to evaluating on moving object segmentation, we also benchmark the model on camouflaged object detection on the MoCA dataset, where visual cues are often less effective than motion cues. To compare fairly with CIS [86], we use the code and model released by the authors, and fine-tune their model on MoCA in a self-supervised manner. We convert our model’s output segmentation mask to a bounding box by keeping only the largest connected region in the prediction before taking the bounding box around it.

We report quantitative results in Table 3 and show qualitative results in Figure 4. Our model significantly outperforms CIS (14% when allowing no post-processing), pre-

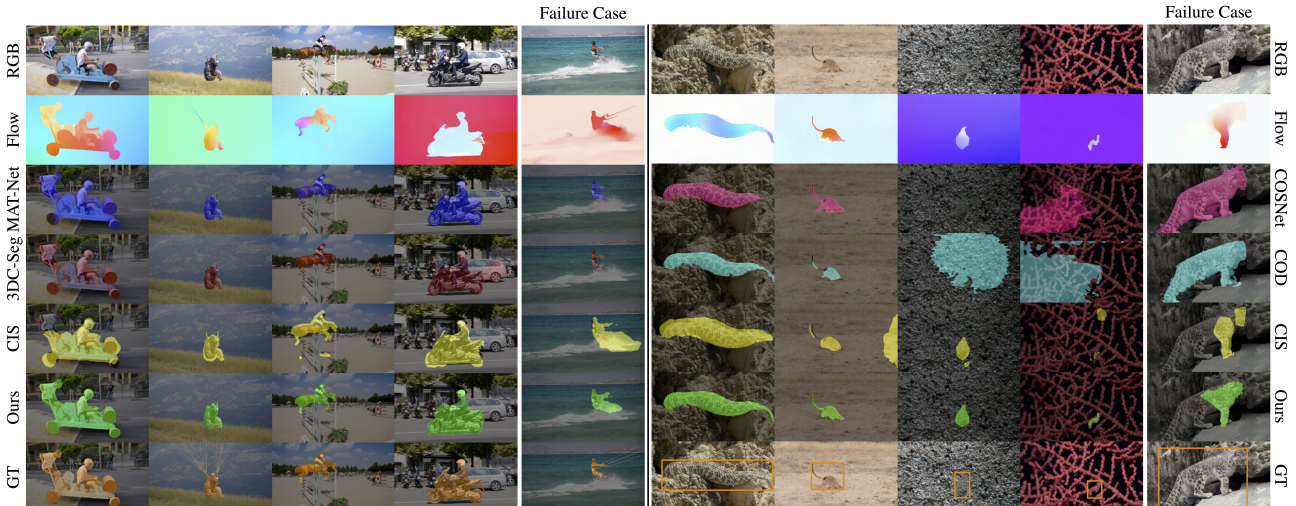


Figure 4. **Qualitative results.** On DAVIS2016 (left), our method is able to segment a variety of challenging objects such as people, animals, and vehicles, often on-par with top supervised approaches. On MoCA (right), our model is able to accurately segment well-camouflaged objects even when previous supervised methods fail completely (3rd, 4th columns). We show a failure case (left) where the splash created by the person is incorrectly included in our predicted segment, and another failure case (right) where the animal is only partially moving and thus partially segmented.

Model	Sup.	RGB	Flow	$\mathcal{J} \uparrow$	Success Rate					SR_{mean}
					$\tau = 0.5$	$\tau = 0.6$	$\tau = 0.7$	$\tau = 0.8$	$\tau = 0.9$	
COD [41]	✓	✗	✓	44.9	0.414	0.330	0.235	0.140	0.059	0.236
COD (two-stream) [41]	✓	✓	✓	55.3	0.602	0.523	0.413	0.267	0.088	0.379
COSNet [48]	✓	✓	✗	50.7	0.588	0.534	0.457	0.337	0.167	0.417
MATNet [89]	✓	✓	✓	64.2	0.712	0.670	0.599	0.492	0.246	0.544
CIS	✗	✓	✓	49.4	0.556	0.463	0.329	0.176	0.030	0.311
CIS (post-processing)	✗	✓	✓	54.1	0.631	0.542	0.399	0.210	0.033	0.363
Ours	✗	✗	✓	63.4	0.742	0.654	0.524	0.351	0.147	0.484

Table 3. **Comparison results on MoCA dataset.** We report the successful localization rate for various thresholds τ (see Section 4.2). Both CIS and Ours were pre-trained on DAVIS and finetuned on MoCA in a self-supervised manner. Note that, our method achieves comparable Jaccard (\mathcal{J}) to MATNet (2nd best model on DAVIS), without using RGB inputs and without any manual annotation for training.

vious supervised approaches *e.g.* COD [41] (18.5% on Jaccard), and even COSNet [48] (among the top supervised approaches on DAVIS). We conjecture that COSNet’s weaker performance is due to its sole reliance on visual appearance (which is distracting for the MoCA data) rather than using motion inputs. This is particularly interesting, as it clearly indicates that no single information cue is able to do the task perfectly, echoing the two-stream hypothesis [28] that both appearance and motion are essential to visual systems.

5.4. Limitations

Despite showing remarkable improvements on motion segmentation in accuracy and runtime, we note the following limitations of the proposed approach (shown in Figure 4) and treat them as future work: *first*, the existing benchmarks are mostly limited to motion segmentation into foreground and background, thus, we choose to use two

slots in this paper; however, in real scenarios, videos may contain multiple independently moving objects, which the current model will assign to a single layer. It may be desirable to further separate these objects into different layers. *Second*, we have only explored motion-only (optical flow) as input, which significantly limits the model in segmenting objects when flow is uninformative or incomplete (as in Fig. 4, right); however, the self-supervised video object segmentation objective is applicable also to a two-stream approach, and so RGB could be incorporated. *Third*, the current method may fail when optical flow is noisy or low-quality (Fig. 4, left); jointly optimizing flow and segmentation is a possible way forward in this case.

6. Conclusion

In this paper, we present a self-supervised model for motion segmentation. The algorithm takes only flow as

input, and is trained without any manual annotation, surpassing previous self-supervised methods on public benchmarks such as DAVIS2016, narrowing the gap with supervised methods. On the more challenging camouflage dataset (MoCA), our model actually compares favourably to the top approaches in video object segmentation that are trained with heavy supervision. As computation power grows and more high quality videos become available, we believe that self-supervised learning algorithms can serve as a strong competitor to the supervised counterparts for their scalability and generalizability.

References

- [1] Jean-Baptiste Alayrac, João Carreira, and Andrew Zisserman. The visual centrifuge: Model-free layered video representations. In *Proc. CVPR*, 2019.
- [2] Jean-Baptiste Alayrac, Joao Carreira, Relja Arandjelovic, and Andrew Zisserman. Controllable attention for structured layered video decomposition. In *Proc. ICCV*, 2019.
- [3] Niki Parmar Jakob Uszkoreit Llion Jones Aidan N. Gomez Lukasz Kaiser Illia Polosukhin Ashish Vaswani, Noam Shazeer. Attention is all you need. In *NIPS*, 2017.
- [4] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? *arXiv preprint arXiv:2102.05095*, 2021.
- [5] Pia Bideau and Erik Learned-Miller. A detailed rubric for motion segmentation. *arXiv preprint arXiv:1610.10033*, 2016.
- [6] Pia Bideau and Erik Learned-Miller. It’s moving! a probabilistic model for causal motion segmentation in moving camera videos. In *Proc. ECCV*, 2016.
- [7] Pia Bideau, Rakesh R Menon, and Erik Learned-Miller. Moa-net: self-supervised motion segmentation. In *Proc. ECCV Workshop*, 2018.
- [8] Gabriel J Brostow and Irfan A Essa. Motion based decompositing of video. In *Proc. ICCV*, 1999.
- [9] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, 2020.
- [10] Thomas Brox and Jitendra Malik. Object segmentation by long term analysis of point trajectories. In *Proc. ECCV*, 2010.
- [11] Christopher P. Burgess, Loic Matthey, Nicholas Watters, Rishabh Kabra, Irina Higgins, Matt Botvinick, and Alexander Lerchner. Monet: Unsupervised scene decomposition and representation. *arXiv preprint arXiv:1901.11390*, 2019.
- [12] Sergi Caelles, Kevis-Kokitsi Maninis, Jordi Pont-Tuset, Laura Leal-Taixé, Daniel Cremers, and Luc Van Gool. One-shot video object segmentation. In *Proc. CVPR*, 2017.
- [13] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Proc. ECCV*, 2020.
- [14] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pre-training from pixels. In *Proceedings of the International Conference on Machine Learning*, 2020.
- [15] J. Cheng, Y.-H. Tsai, S. Wang, and M.-H. Yang. Segflow: Joint learning for video object segmentation and optical flow. In *Proc. ICCV*, 2017.
- [16] Achal Dave, Pavel Tokmakov, and Deva Ramanan. Towards segmenting anything that moves. In *Proc. ICCV Workshop*, 2019.
- [17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. Imagenet: A large-scale hierarchical image database. In *Proc. CVPR*, 2009.
- [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [19] Carl Doersch and Andrew Zisserman. Sim2real transfer learning for 3d human pose estimation: motion to the rescue. *Proc. NeurIPS*, 2019.
- [20] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *Proc. ICLR*, 2021.
- [21] Martin Engelcke, Adam R Kosiorek, Oiwi Parker Jones, , and Ingmar Posner. Genesis: Generative scene inference and sampling with object-centric latent representations. In *Proc. ICLR*, 2020.
- [22] Mark Everingham, Luc Van Gool, Chris K. I. Williams, John Winn, and Andrew Zisserman. The PASCAL Visual Object Classes (VOC) challenge. *International Journal of Computer Vision*, 88(2):303–338, 2010.
- [23] Alon Faktor and Michal Irani. Video segmentation by non-local consensus voting. In *Proc. BMVC*, 2014.
- [24] Deng-Ping Fan, Wenguan Wang, Ming-Ming Cheng, and Jianbing Shen. Shifting more attention to video salient object detection. In *Proc. CVPR*, 2019.
- [25] Katerina Fragkiadaki, Geng Zhang, and Jianbo Shi. Video segmentation by tracing discontinuities in a trajectory embedding. In *Proc. CVPR*, 2012.
- [26] Yossi Gandelsman, Assaf Shocher, and Michal Irani. “Double-DIP”: Unsupervised image decomposition via coupled deep-image-priors. In *Proc. CVPR*, 2019.
- [27] Rohit Girdhar, João Carreira, Carl Doersch, and Andrew Zisserman. Video Action Transformer Network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [28] Melvyn A. Goodale and A. David Milner. Separate visual pathways for perception and action. *Trends in Neurosciences*, 15(1):20–25, 1992.
- [29] Klaus Greff, Raphael Lopez Kaufman, Rishabh Kabra, Nick Watters, Chris Burgess, Daniel Zoran, Loic Matthey, Matthew Botvinick, and Alexander Lerchner. Multi-object representation learning with iterative variational inference.

- In *Proc. ICML*, 2019.
- [30] Suyog Jain, Bo Xiong, and Kristen Grauman. Fusionseg: Learning to combine motion and appearance for fully automatic segmentation of generic objects in videos. *arXiv preprint arXiv:1701.05384*, 2017.
- [31] Suyog Dutt Jain, Bo Xiong, and Kristen Grauman. Fusionseg: Learning to combine motion and appearance for fully automatic segmentation of generic objects in videos. In *Proc. CVPR*, 2017.
- [32] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proc. CVPR*, 2017.
- [33] Njegica Jojic and B.J. Frey. Learning flexible sprites in video layers. In *Proc. CVPR*, 2001.
- [34] Yeong Jun Koh and Chang-Su Kim. Primary object segmentation in videos based on region augmentation and reduction. In *Proc. CVPR*, 2017.
- [35] Margret Keuper, Bjoern Andres, and Thomas Brox. Motion trajectory segmentation via minimum cost multicuts. In *Proc. ICCV*, 2015.
- [36] Margret Keuper, Bjoern Andres, and Thomas Brox. Motion trajectory segmentation via minimum cost multicuts. In *Proc. ICCV*, 2015.
- [37] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- [38] M Pawan Kumar, Philip HS Torr, and Andrew Zisserman. Learning layered motion segmentations of video. *International Journal of Computer Vision*, 76(3):301–319, 2008.
- [39] Zihang Lai, Erika Lu, and Weidi Xie. MAST: A memory-augmented self-supervised tracker. In *Proc. CVPR*, 2020.
- [40] Zihang Lai and Weidi Xie. Self-supervised learning for video correspondence flow. In *BMVC*, 2019.
- [41] Hala Lamdouar, Charig Yang, Weidi Xie, and Andrew Zisserman. Betrayed by motion: Camouflaged object discovery via motion segmentation. *Proc. ACCV*, 2020.
- [42] Fuxin Li, Taeyoung Kim, Ahmad Humayun, David Tsai, and James M. Rehg. Video segmentation by tracking many figure-ground segments. In *Proc. ICCV*, 2013.
- [43] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollar. Microsoft coco: Common objects in context. In *Proceedings of the European Conference on Computer Vision*, 2014.
- [44] Liang Liu, Jiangning Zhang, Ruifei He, Yong Liu, Yabiao Wang, Ying Tai, Donghao Luo, Chengjie Wang, Jilin Li, and Feiyue Huang. Learning by analogy: Reliable supervision from transformations for unsupervised optical flow estimation. In *Proc. CVPR*, 2020.
- [45] Pengpeng Liu, Michael R. Lyu, Irwin King, and Jia Xu. Self-low: Self-supervised learning of optical flow. In *Proc. CVPR*, 2019.
- [46] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. In *Proc. NeurIPS*, 2020.
- [47] Erika Lu, Forrester Cole, Tali Dekel, Weidi Xie, Andrew Zisserman, David Salesin, William T Freeman, and Michael Rubinstein. Layered neural rendering for retiming people in video. In *SIGGRAPH Asia*, 2020.
- [48] Xiankai Lu, Wenguan Wang, Chao Ma, Jianbing Shen, Ling Shao, and Fatih Porikli. See more, know more: Unsupervised video object segmentation with co-attention siamese networks. In *Proc. CVPR*, 2019.
- [49] Sabarinath Mahadevan, Ali Athar, Aljoša Ošep, Sebastian Hennen, Laura Leal-Taixé, and Bastian Leibe. Making a case for 3d convolutions for object segmentation in videos. In *Proc. BMVC*, 2020.
- [50] Aravindh Mahendran, James Thewlis, and Andrea Vedaldi. Cross pixel optical-flow similarity for self-supervised learning. In *Proceedings of the Asian Conference on Computer Vision*, 2018.
- [51] K-K Maninis, Sergi Caelles, Yuhua Chen, Jordi Pont-Tuset, Laura Leal-Taixé, Daniel Cremers, and Luc Van Gool. Video object segmentation without temporal information. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2018.
- [52] Peter Ochs and Thomas Brox. Object segmentation in video: a hierarchical variational approach for turning point trajectories into dense regions. In *Proc. ICCV*, 2011.
- [53] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *Proc. ICCV*, 2019.
- [54] Anestis Papazoglou and Vittorio Ferrari. Fast object segmentation in unconstrained video. In *Proc. ICCV*, 2013.
- [55] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Proc. CVPR*, 2016.
- [56] Ochs Peter, Malik Jitendra, and Brox Thomas. Segmentation of moving objects by long term video analysis. *TPAMI*, 2014.
- [57] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alex Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv preprint arXiv:1704.00675*, 2017.
- [58] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.
- [59] VS Ramachandran. Guest editorial: The neurobiology of perception. *Perception*, 14(97):103, 1985.
- [60] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. *arXiv preprint arXiv:2102.12092*, 2021.
- [61] Anurag Ranjan, Varun Jampani, Lukas Balles, Deqing Sun, Kihwan Kim, Jonas Wulff, and Michael J. Black. Competitive collaboration: Joint unsupervised learning of depth, camera motion, optical flow and motion segmentation. In *Proc. CVPR*, 2019.
- [62] Hedvig Sidenbladh, Michael J Black, and David J Fleet. Stochastic tracking of 3d human figures using 2d image motion. In *Proc. ECCV*, 2000.
- [63] Josef Sivic, Frederik Schaffalitzky, and Andrew Zisserman. Object level grouping for video shots. In *Proceedings of the 8th European Conference on Computer Vision, Prague, Czech Republic*. Springer-Verlag, 2004.

- [64] Hongmei Song, Wenguan Wang, Sanyuan Zhao, Jianbing Shen, and Kin-Man Lam. Pyramid dilated deeper convlstm for video salient object detection. In *Proc. ECCV*, 2018.
- [65] Pratul P. Srinivasan, Richard Tucker, Jonathan T. Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely. Pushing the boundaries of view extrapolation with multiplane images. In *Proc. CVPR*, 2019.
- [66] D. Sun, E. Sudderth, and M. J. Black. Layered segmentation and optical flow estimation over time. In *Proc. CVPR*, 2012.
- [67] Deqing Sun, Jonas Wulff, Erik Sudderth, Hanspeter Pfister, and Michael Black. A fully-connected layered model of foreground and background flow. In *Proc. CVPR*, 2013.
- [68] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume. In *Proc. CVPR*, 2018.
- [69] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Proc. ECCV*, 2020.
- [70] Pavel Tokmakov, Karteek Alahari, and Cordelia Schmid. Learning motion patterns in videos. In *Proc. CVPR*, 2017.
- [71] Pavel Tokmakov, Cordelia Schmid, and Karteek Alahari. Learning to segment moving objects. *International Journal of Computer Vision*, 2019.
- [72] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- [73] Paul Voigtlaender, Yuning Chai, Florian Schroff, Hartwig Adam, Bastian Leibe, and Liang-Chieh Chen. Feelvos: Fast end-to-end embedding learning for video object segmentation. In *Proc. CVPR*, 2019.
- [74] Paul Voigtlaender and Bastian Leibe. Online adaptation of convolutional neural networks for video object segmentation. 2017.
- [75] Carl Vondrick, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy. Tracking emerges by colorizing videos. In *Proc. ECCV*, 2018.
- [76] J. Y. A. Wang and E. H. Adelson. Representing moving images with layers. *The IEEE Transactions on Image Processing Special Issue: Image Sequence Compression*, 3(5):625–638, September 1994.
- [77] Wenguan Wang, Xiankai Lu, Jianbing Shen, David J Crandall, and Ling Shao. Zero-shot video object segmentation via attentive graph neural networks. In *Proc. ICCV*, 2019.
- [78] Wenguan Wang, Jianbing Shen, Ruigang Yang, and Fatih Porikli. Saliency-aware video object segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(1):20–33, 2018.
- [79] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proc. CVPR*, 2018.
- [80] Max Wertheimer. Untersuchungen zur lehre von der gestalt. ii. *Psychologische forschung*, 4(1):301–350, 1923.
- [81] Jonas Wulff and Michael J. Black. Efficient sparse-to-dense optical flow estimation using a learned basis and layers. In *Proc. ECCV*.
- [82] Jonas Wulff and Michael J. Black. Modeling blurred video with layers. In *Proc. ICCV*, 2014.
- [83] Christopher Xie, Yu Xiang, Zaid Harchaoui, and Dieter Fox. Object discovery in videos as foreground motion clustering. In *Proc. CVPR*, 2019.
- [84] Ning Xu, Linjie Yang, Yuchen Fan, Dingcheng Yue, Yuchen Liang, Jianchao Yang, and Thomas Huang. Youtube-vos: A large-scale video object segmentation benchmark. In *Proc. ECCV*, 2018.
- [85] Tianfan Xue, Michael Rubinstein, Ce Liu, and William T. Freeman. A computational approach for obstruction-free photography. *ACM Transactions on Graphics (Proc. SIGGRAPH)*, 34(4), 2015.
- [86] Yanchao Yang, Antonio Loquercio, Davide Scaramuzza, and Stefano Soatto. Unsupervised moving object detection via contextual information separation. In *Proc. CVPR*.
- [87] Zhao Yang, Qiang Wang, Luca Bertinetto, Song Bai, Weiming Hu, and Philip H.S. Torr. Anchor diffusion for unsupervised video object segmentation. In *Proc. ICCV*, 2019.
- [88] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *SIGGRAPH*, 2018.
- [89] Tianfei Zhou, Shunzhou Wang, Yi Zhou, Yazhou Yao, Jianwu Li, and Ling Shao. Motion-attentive transition for zero-shot video object segmentation. In *AAAI*, volume 34, pages 13066–13073, 2020.
- [90] C Lawrence Zitnick, Sing Bing Kang, Matthew Uyttendaele, Simon Winder, and Richard Szeliski. High-quality video view interpolation using a layered representation. *ACM transactions on graphics (TOG)*, 23(3):600–608, 2004.

A. Training Details

In this section, we outline the training settings and architecture details used in our paper. Codes and models will be released.

A.1. Encoder & Decoder

The backbone of network architecture for the model are shown in Table 4, we refer the readers to the pseudocode for iterative binding module.

	stage	operation	output sizes
	input	–	$3 \times 128 \times 224$
Encoder	conv1	$[5 \times 5, 64] \times 2$	$64 \times 128 \times 224$
	mp1	maxpool, stride = 2	$64 \times 64 \times 112$
	conv2	$[5 \times 5, 128] \times 2$	$128 \times 64 \times 112$
	mp2	maxpool, stride = 2	$128 \times 32 \times 56$
	conv3	$[5 \times 5, 256] \times 2$	$256 \times 32 \times 56$
Decoder	conv ^T 1	$5 \times 5, 64, \text{stride} = 2$	$64 \times 2 \times 32 \times 56$
	conv ^T 2	$5 \times 5, 64, \text{stride} = 2$	$64 \times 2 \times 64 \times 112$
	conv ^T 3	$5 \times 5, 64, \text{stride} = 2$	$64 \times 2 \times 128 \times 224$
	outconv	$5 \times 5, 64,$	$4 \times 2 \times 128 \times 224$
		$5 \times 5, 4$	

Table 4. Network architecture. All convolutions have padding 2 to preserve spatial resolution, and are followed by instance normalization and ReLU activation, except the final layer. The details of the iterative binding module is in Figure 5.

A.2. Iterative Binding Module

The pseudocode for the iterative binding module is shown in Figure 5. The full code is available in the submitted source code, and will be made publicly available.

Algorithm: Pseudo Code for Iterative Binding / Grouping
<code># inputs: feature maps + position encoding</code>
<code>slots = Embedding(D, 2) # [bsz, 2, D]</code>
<code>K = project_K(inputs) # [bsz, HW, D]</code>
<code>V = project_V(inputs) # [bsz, HW, D]</code>
<code>for _ in range(T): # T iterations</code>
<code> slots = LayerNorm(slots) # [bsz, 2, D]</code>
<code> Q = project_Q(slots) # [bsz, 2, D]</code>
<code> scores = dot(K, Q.t()) # [bsz, HW, 2]</code>
<code> attn = softmax(scores / sqrt(D), axis=-1) # [bsz, HW, 2]</code>
<code> updates = weighted_mean(attn.t(), V) # [bsz, 2, D]</code>
<code> slots = GRU(slots, updates) # [bsz, 2, D]</code>
<code> slots = slots + MLP(LayerNorm(slots)) # [bsz, 2, D]</code>

Figure 5. Pseudocode for iterative binding. The linear projections for key (K) and value (V) have 256 dimensions. The MLP has two layers, both with 256 dimensions, with ReLU in between.

A.3. Hyperparameters

For all datasets, we train with batch size of 64, although note that this corresponds to 32 pairs of optical flow in order to train with consistency loss. Our initial learning rate is $5e-4$, with first 200 steps being warmup, and decays by half every $8e4$ iterations. During this decay, the scale for entropy and consistency loss is also increased by a factor of 5, gradually encouraging the predicted alpha channel to be binary. We train the algorithm for about 300k iterations.

B. MoCA dataset curation

The MoCA dataset contains 141 high-definition video sequences, with an average duration of 11 seconds. These sequences were collected from YouTube with resolution 720×1280 , and sampled at 24 fps, resulting in 37K frames depicting 67 kinds of camouflaged animals moving in natural scenes. Both temporal and spatial annotations are provided in the form of tight bounding boxes on every 5th frame. We use a modified version of this dataset in order to make it more suitable for segmentation tasks. We outline the process and rationale below.

- We crop away the channel logos and empty border spaces, and resize the low-resolution images to the same resolution as the other images in the dataset (at 720×1280). We then adjust the ground-truth annotations accordingly.
- The original authors resampled all videos to 24 fps even when some original videos have less, causing some consecutive frames to be identical. To alleviate this, we sample every 3 frames from the original dataset, up to 100 frames per video.
- For annotations, we use linear interpolation to generate the missing frames' bounding boxes, resulting in a dense frame-wise annotation.
- The authors also provided the motion labels for each annotated frame (locomotion, deformation, static), so we filter away videos with predominantly no locomotion.
- We also further discard videos that contain large amount of frames where the motion does not belong to the primary object.

This eventually results in 88 video sequences and 4803 frames, which we will release for fair comparison.

C. Results breakdown

The main evaluation metric used in this paper is the Jaccard score (\mathcal{J}), which is the intersection-over-union between the predicted and ground-truth masks. In line with previous works, we show the per-category results breakdown for our model on DAVIS2016, SegTrackv2, FBMS59, and MoCA in Tables 5-8. Note that, since we focus on both speed and accuracy, the predictions are only of 128×224 pixels, and we directly upsample this prediction to the original resolution and compare with the groundtruth.

Sequence	$\mathcal{J}(M)$	$\mathcal{J}(R)$	$\mathcal{J}(D)$	$\mathcal{F}(M)$	$\mathcal{F}(R)$	$\mathcal{F}(D)$
bear	0.766	1.000	-4.5	0.640	0.914	-1.9
blackswan	0.795	1.000	5.9	0.658	0.980	7.3
bm-x-bumps	0.373	0.114	8.3	0.260	0.011	10.8
bm-x-trees	0.744	1.000	14.5	0.666	0.909	13.9
boat	0.602	0.787	44.5	0.724	0.936	22.4
breakdance	0.792	0.990	-2.4	0.814	1.000	10.8
breakdance-flare	0.826	1.000	10.2	0.771	1.000	16.7
bus	0.391	0.000	-4.4	0.248	0.000	-5.0
camel	0.540	0.679	9.3	0.667	0.923	-9.8
car-roundabout	0.479	0.500	1.2	0.404	0.146	-10.6
car-shadow	0.879	1.000	-0.2	0.818	1.000	-3.3
car-turn	0.492	0.537	3.4	0.560	0.683	6.8
cows	0.690	0.898	-5.0	0.653	0.852	4.4
dance-jump	0.815	1.000	-4.8	0.683	1.000	5.9
dance-twirl	0.739	0.960	7.4	0.803	0.900	28.5
dog	0.701	0.789	2.9	0.477	0.526	9.3
dog-agility	0.810	1.000	7.7	0.858	1.000	4.0
drift-chicane	0.777	1.000	7.3	0.611	0.646	26.3
drift-straight	0.864	1.000	-5.2	0.685	1.000	-6.8
drift-turn	0.593	0.654	27.7	0.220	0.000	20.4
Average	0.683	0.795	6.2	0.611	0.721	7.5

Table 5. Full results on DAVIS2016. J refers to the Jaccard score, while the F-measure refers to the contour accuracy. M, R, and D refers to mean, recall and decay respectively.

Sequence	$\mathcal{J}(M)$
bird of paradise	0.791
birdfall	0.300
bm-x	0.609
cheetah	0.370
drift	0.797
frog	0.733
girl	0.746
hummingbird	0.506
monkey	0.751
monkeydog	0.133
parachute	0.914
penguin	0.697
soldier	0.741
worm	0.326
seq avg	0.601
frames avg	0.586

Table 6. Sequence-wise results on SegTrackv2.

Sequence	$\mathcal{J}(M)$
camel01	0.281
cars1	0.846
cars10	0.322
cars4	0.826
cars5	0.842
cats01	0.672
cats03	0.640
cats06	0.362
dogs01	0.629
dogs02	0.636
farm01	0.816
giraffes01	0.322
goats01	0.375
horses02	0.628
horses04	0.566
horses05	0.334
lion01	0.399
marple12	0.680
marple2	0.750
marple4	0.799
marple6	0.450
marple7	0.567
marple9	0.537
people03	0.598
people1	0.761
people2	0.842
rabbits02	0.415
rabbits03	0.319
rabbits04	0.400
tennis	0.561
seq avg	0.573
frames avg	0.531

Table 7. Sequence-wise results on FBMS59.

D. Qualitative results

We show more qualitative results in Figures 6 and 7. We also submit video sequences of the results as part of the Supplementary Material.

sequence	\mathcal{J}	$\tau_{0.5}$	$\tau_{0.6}$	$\tau_{0.7}$	$\tau_{0.8}$	$\tau_{0.9}$	<i>avg</i>
arabian_horn_viper	0.709	0.990	0.909	0.586	0.091	0.000	0.515
arctic_fox	0.381	0.404	0.383	0.298	0.191	0.000	0.255
arctic_fox_1	0.879	0.913	0.913	0.913	0.913	0.652	0.861
arctic_wolf_0	0.705	0.929	0.859	0.596	0.202	0.051	0.527
arctic_wolf_1	0.712	0.795	0.795	0.795	0.744	0.256	0.677
bear	0.508	0.611	0.453	0.221	0.074	0.000	0.272
black_cat_0	0.499	0.556	0.476	0.302	0.048	0.000	0.276
black_cat_1	0.086	0.030	0.020	0.010	0.000	0.000	0.012
crab	0.594	0.800	0.400	0.200	0.000	0.000	0.280
crab_1	0.288	0.309	0.200	0.073	0.000	0.000	0.116
cuttlefish_0	0.222	0.194	0.032	0.000	0.000	0.000	0.045
cuttlefish_1	0.034	0.043	0.000	0.000	0.000	0.000	0.009
cuttlefish_4	0.655	1.000	0.846	0.231	0.000	0.000	0.415
cuttlefish_5	0.724	0.870	0.870	0.783	0.304	0.043	0.574
dead_leaf_butterfly_1	0.784	0.913	0.783	0.739	0.696	0.304	0.687
desert_fox	0.470	0.362	0.234	0.191	0.149	0.106	0.209
devil_scorpionfish	0.938	1.000	1.000	1.000	0.913	0.826	0.948
devil_scorpionfish_1	0.913	1.000	1.000	1.000	1.000	0.565	0.913
devil_scorpionfish_2	0.857	0.968	0.903	0.871	0.806	0.548	0.819
egyptian_nightjar	0.765	0.905	0.842	0.789	0.611	0.158	0.661
elephant	0.728	0.783	0.652	0.609	0.565	0.348	0.591
flatfish_0	0.575	0.677	0.657	0.636	0.485	0.141	0.519
flatfish_1	0.682	0.848	0.835	0.722	0.354	0.051	0.562
flatfish_2	0.765	0.839	0.774	0.774	0.742	0.645	0.755
flatfish_4	0.697	0.958	0.895	0.568	0.126	0.000	0.509
flounder	0.896	1.000	1.000	0.986	0.986	0.437	0.882
flounder_3	0.505	0.429	0.286	0.143	0.143	0.000	0.200
flounder_4	0.767	0.949	0.897	0.769	0.487	0.128	0.646
flounder_5	0.681	0.797	0.747	0.696	0.468	0.165	0.575
flounder_6	0.683	0.768	0.677	0.616	0.465	0.263	0.558
flounder_7	0.719	0.930	0.845	0.662	0.254	0.028	0.544
flounder_8	0.707	0.925	0.774	0.645	0.409	0.000	0.551
flounder_9	0.636	0.821	0.718	0.436	0.231	0.026	0.446
fossa	0.280	0.143	0.000	0.000	0.000	0.000	0.029
goat_0	0.589	0.707	0.636	0.414	0.212	0.020	0.398
goat_1	0.744	0.930	0.930	0.704	0.380	0.085	0.606
groundhog	0.525	0.646	0.525	0.374	0.162	0.010	0.343
hedgehog_0	0.329	0.298	0.170	0.085	0.021	0.021	0.119
hedgehog_1	0.471	0.533	0.400	0.333	0.067	0.067	0.280
hedgehog_2	0.771	0.800	0.733	0.733	0.600	0.267	0.627
hedgehog_3	0.486	0.564	0.359	0.282	0.128	0.026	0.272
hermit_crab	0.674	0.806	0.806	0.677	0.419	0.065	0.555
ibex	0.390	0.513	0.359	0.128	0.000	0.000	0.200
jerboa	0.555	0.739	0.435	0.348	0.130	0.000	0.330

sequence	\mathcal{J}	$\tau_{0.5}$	$\tau_{0.6}$	$\tau_{0.7}$	$\tau_{0.8}$	$\tau_{0.9}$	<i>avg</i>
jerboa_1	0.450	0.452	0.323	0.226	0.097	0.000	0.219
lichen_katydid	0.502	0.455	0.323	0.162	0.020	0.010	0.194
lion_cub_0	0.728	0.972	0.873	0.549	0.282	0.127	0.561
lion_cub_1	0.571	0.687	0.556	0.354	0.141	0.010	0.349
lion_cub_3	0.179	0.182	0.141	0.081	0.051	0.000	0.091
lioness	0.423	0.419	0.129	0.000	0.000	0.000	0.110
marine_iguana	0.376	0.217	0.130	0.000	0.000	0.000	0.070
markhor	0.808	0.909	0.855	0.800	0.709	0.364	0.727
meerkat	0.778	0.871	0.774	0.742	0.710	0.323	0.684
mountain_goat	0.742	1.000	0.968	0.710	0.226	0.065	0.594
nile_monitor_1	0.524	0.616	0.434	0.283	0.121	0.000	0.291
octopus	0.637	0.838	0.687	0.273	0.131	0.020	0.390
octopus_1	0.411	0.242	0.091	0.061	0.061	0.010	0.093
peacock_flounder_0	0.884	0.931	0.931	0.931	0.931	0.701	0.885
peacock_flounder_1	0.812	0.970	0.929	0.859	0.667	0.222	0.729
peacock_flounder_2	0.873	1.000	1.000	0.989	0.926	0.368	0.857
polar_bear_0	0.622	0.845	0.676	0.352	0.141	0.000	0.403
polar_bear_1	0.487	0.603	0.556	0.476	0.175	0.016	0.365
polar_bear_2	0.792	0.949	0.846	0.744	0.641	0.308	0.697
pygmy_seahorse_2	0.478	0.582	0.364	0.255	0.036	0.000	0.247
pygmy_seahorse_4	0.654	0.935	0.839	0.516	0.000	0.000	0.458
rodent_x	0.738	0.870	0.826	0.696	0.435	0.174	0.600
scorpionfish_0	0.622	0.761	0.648	0.521	0.408	0.127	0.493
scorpionfish_1	0.616	0.681	0.574	0.426	0.319	0.128	0.426
scorpionfish_2	0.842	0.962	0.949	0.924	0.823	0.380	0.808
scorpionfish_3	0.804	0.975	0.861	0.785	0.595	0.354	0.714
scorpionfish_4	0.800	1.000	0.949	0.821	0.513	0.231	0.703
scorpionfish_5	0.899	1.000	1.000	1.000	0.957	0.478	0.887
seal_1	0.865	1.000	0.913	0.870	0.870	0.652	0.861
seal_2	0.676	0.800	0.655	0.455	0.291	0.145	0.469
seal_3	0.445	0.400	0.200	0.067	0.000	0.000	0.133
shrimp	0.772	0.933	0.867	0.733	0.667	0.133	0.667
snow_leopard_0	0.760	1.000	0.949	0.692	0.359	0.077	0.615
snow_leopard_1	0.772	0.826	0.696	0.652	0.652	0.522	0.670
snow_leopard_2	0.883	1.000	0.989	0.968	0.863	0.526	0.869
snow_leopard_3	0.608	0.778	0.556	0.417	0.333	0.083	0.433
snow_leopard_6	0.816	0.830	0.787	0.787	0.723	0.660	0.757
snow_leopard_7	0.679	0.871	0.806	0.581	0.258	0.000	0.503
snow_leopard_8	0.556	0.702	0.596	0.447	0.191	0.000	0.387
snowy_owl_0	0.564	0.532	0.340	0.298	0.128	0.000	0.260
spider_tailed_horned_viper_0	0.473	0.400	0.333	0.267	0.067	0.067	0.227
spider_tailed_horned_viper_1	0.558	0.620	0.423	0.310	0.254	0.056	0.332
spider_tailed_horned_viper_2	0.848	0.964	0.945	0.909	0.782	0.564	0.833
spider_tailed_horned_viper_3	0.860	1.000	1.000	1.000	0.677	0.387	0.813
seq avg	0.634	0.734	0.640	0.522	0.361	0.166	0.485
frames avg	0.634	0.742	0.654	0.524	0.351	0.147	0.484

Table 8. Results breakdown for MoCA.

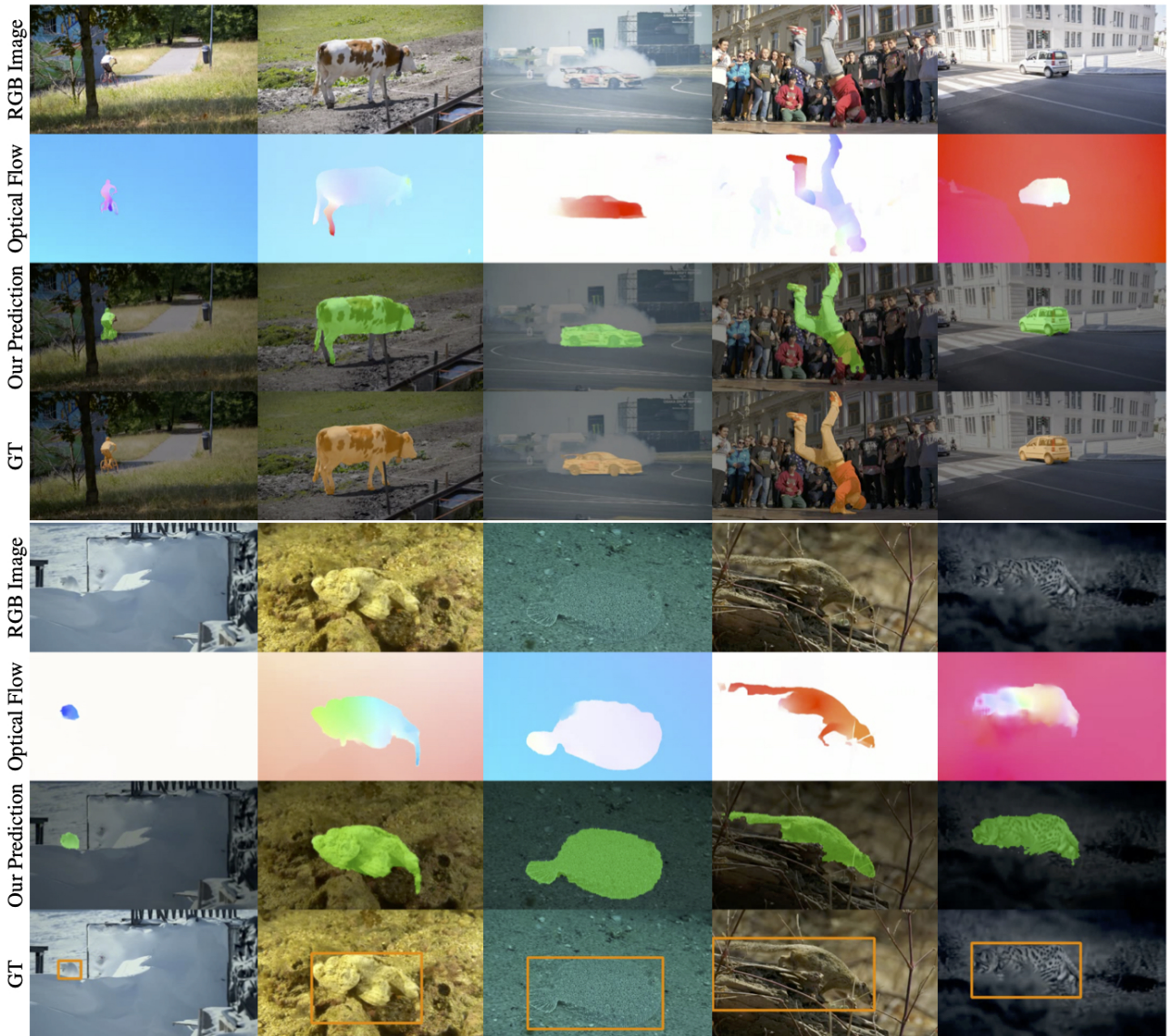


Figure 6. Qualitative results on DAVIS2016 and MoCA, respectively.



Figure 7. Qualitative results on SegTrackv2 and FBMS59, respectively.


Statement of Authorship for joint/multi-authored papers for PGR thesis

To appear at the end of each thesis chapter submitted as an article/paper

The statement shall describe the candidate's and co-authors' independent research contributions in the thesis publications. For each publication there should exist a complete statement that is to be filled out and signed by the candidate and supervisor (**only required where there isn't already a statement of contribution within the paper itself**).

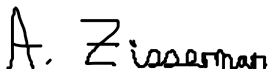
Title of Paper	Self-supervised Video Object Segmentation by Motion Grouping
Publication Status	<input type="checkbox"/> Published <input type="checkbox"/> Accepted for Publication <input type="checkbox"/> Submitted for Publication <input type="checkbox"/> Unpublished and unsubmitted work written in a manuscript style
Publication Details	Yang, C., Lamdouar, H., Lu, E. , Zisserman, A., & Xie, W. In Proceedings of the IEEE International Conference on Computer Vision (ICCV), 2021.

Student Confirmation

Student Name:	Erika Lu		
Contribution to the Paper	<ul style="list-style-type: none">• Development of idea• Writing the paper		
Signature		Date	25 th August, 2021

Supervisor Confirmation

By signing the Statement of Authorship, you are certifying that the candidate made a substantial contribution to the publication, and that the description described above is accurate.

Supervisor name and title:	Professor Andrew Zisserman		
Supervisor comments			
Signature		Date	9/9/2021

This completed form should be included in the thesis, at the end of the relevant chapter.

7

Class-Agnostic Counting

Here we take a brief detour from the topic of video understanding and instead consider the problem of counting objects of *any* class, even those unseen. This work was presented at the Asian Conference on Computer Vision (ACCV), 2018.

Class-Agnostic Counting

Erika Lu, Weidi Xie, and Andrew Zisserman

Visual Geometry Group, University of Oxford
{erika,weidi,az}@robots.ox.ac.uk

Abstract. Nearly all existing counting methods are designed for a specific object class. Our work, however, aims to create a counting model able to count any class of object. To achieve this goal, we formulate counting as a matching problem, enabling us to exploit the image self-similarity property that naturally exists in object counting problems.

We make the following three contributions: *first*, a Generic Matching Network (GMN) architecture that can potentially count any object in a class-agnostic manner; *second*, by reformulating the counting problem as one of matching objects, we can take advantage of the abundance of video data labeled for tracking, which contains natural repetitions suitable for training a counting model. Such data enables us to train the GMN. *Third*, to customize the GMN to different user requirements, an adapter module is used to specialize the model with minimal effort, i.e. using a few labeled examples, and adapting only a small fraction of the trained parameters. This is a form of few-shot learning, which is practical for domains where labels are limited due to requiring expert knowledge (e.g. microbiology).

We demonstrate the flexibility of our method on a diverse set of existing counting benchmarks: specifically cells, cars, and human crowds. The model achieves competitive performance on cell and crowd counting datasets, and surpasses the state-of-the-art on the car dataset using only three training images. When training on the entire dataset, the proposed method outperforms all previous methods by a large margin.

1 Introduction

The objective of this paper is to count objects of interest in an image. In the literature, object counting methods are generally cast into two categories: detection-based counting [5,10,16] or regression-based counting [2,4,8,19,21,24,34]. The former relies on a visual object detector that can localize object instances in an image; this method, however, requires training individual detectors for different objects, and the detection problem remains challenging if only a small number of annotations are given. The latter avoids solving the hard detection problem – instead, methods are designed to learn either a mapping from global image features to a scalar (number of objects), or a mapping from dense image features to a density map, achieving better results on counting overlapping instances. However, previous methods for both categories of method (detection, regression) have only developed algorithms that can count a particular class of objects (e.g. cars, cells, penguins, people).

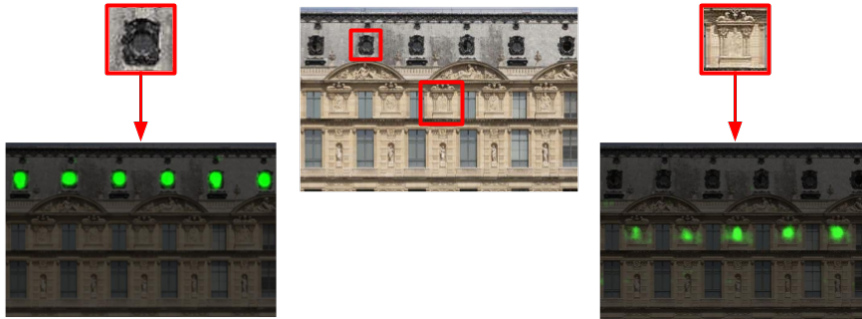


Fig. 1: The model (trained on tracking data) can count an object, e.g. windows or columns, specified as an exemplar patch (in red), without additional training. The heat maps indicate the localizations of the repeated objects. This image is unseen during training.

The objective of this paper is a class-agnostic counting network – one that is able to flexibly count object instances in an image by, for example, simply specifying an exemplar patch of interest as illustrated in Figure 1. To achieve this, we build on a property of images that has been largely ignored explicitly in previous counting approaches – that of *image self-similarity*. At a simplistic level, an image is deemed self-similar if patches repeat to some approximation – for example if patches can be represented by other patches in the same image. Self-similarity has underpinned applications for many vision tasks, ranging from texture synthesis [11], to image denoising [7], to super-resolution [13].

Giving the observation of self-similarity, image counting can be recast as an image *matching* problem – counting instances is performed by matching (self-similar patches) within the same image. To this end we develop a *Generic Matching Network* (GMN) that learns a discriminative classifier to match instances of the exemplar. Furthermore, since matching variations of an object instance within an image is similar to matching variations of an object instance between images, we can take advantage of the abundance of video data labeled for tracking which contains natural repetitions, to train the GMN. This observation, that matching within an image can be thought of as tracking within an image, was previously made by Leung and Malik [22] for the case of repeated elements in an image.

Beyond generic counting, there is often a need to *specialize* matching to more restrictive or general requirements. For example, to count only red cars (rather than all cars) or to count cars at all orientations (which goes beyond simple similarity measures such as squared sum of differences), extending the intra-class variation for the object category of interest [14,18]. To this end, we include an adaptor module that enables fast domain adaptation [28] and few-shot learning [32,33], through the training of a small number of tunable parameters, using very few annotated data.

In the following sections, we begin by detailing the design and training procedure of the GMN in § 2, and demonstrate its capabilities on a set of example counting tasks. In § 3, we adapt the GMN to specialize on several counting benchmark datasets, including the VGG synthetic cells, HeLa cells, and cars captured by drones. During adaptation, only a small number of parameters (3% of the network size) are added and trained on the target domain. Using a very small number of training samples (as few as 3 images for the car dataset), the results achieved are either comparable to, or surpass the current state-of-the-art methods by a large margin. In § 4, we further extend the counting-by-matching idea to a more challenging scenario: Shanghaitech crowd counting, and demonstrate promising results by matching image statistics on scenes where accurate instance-level localization is unobtainable.

2 Method

In this paper, we consider the problem of instance counting, where the objects to be counted in a single query are from the same category, such as the windows on a building, cars in a parking lot, or cells of a certain type.

To exploit the *self-similarity* property, the counting problem is reformalized as localizing and counting “repeated” instances by *matching*. We propose a novel architecture – GMN, and a counting approach which requires learning a comparison scheme for two given objects (patches) in a metric space. The structure of the model naturally accommodates class-agnostic counting, as it learns to search for repetitions of an exemplar patch containing the desired instance. Note that, the concept of *repetition* is defined in a very broad sense; in the following experiments, we show that objects with various shapes, overlaps, and complicated appearance changes can still be treated as “repeated” instances.

The entire GMN consists of three modules, namely, *embedding*, *matching*, and *adapting*, as illustrated in Figure 2. In the *embedding* module, a two-stream network is used to encode the exemplar image patch and the full-resolution image into a feature vector and dense feature map, respectively. In the *matching* module, we learn a discriminative classifier to densely match the exemplar patch to instances in the image. Such learning overcomes within image variations such as illumination changes, small rotations, etc. The object locations and final count can then be acquired by simply taking the *local maximums* or *integration* over the output similarity maps, respectively. Empirically, integral-based counting shows better performance in scenarios where instances have significant overlap, while local max counting is preferred where objects are well-separated, and the positional information is of interest for further processing (e.g. seeds for segmentation).

To avoid the time-consuming collection of annotations for counting data, we use the observation that repetitions occur naturally in videos, as objects are seen under varying viewing conditions from frame to frame. Consequently, we can train the generic matching network with the extensive training data available

for tracking (specifically the ILSVRC video dataset for object detection [31]). In total, the dataset contains nearly 4500 videos and over 1M annotated frames.

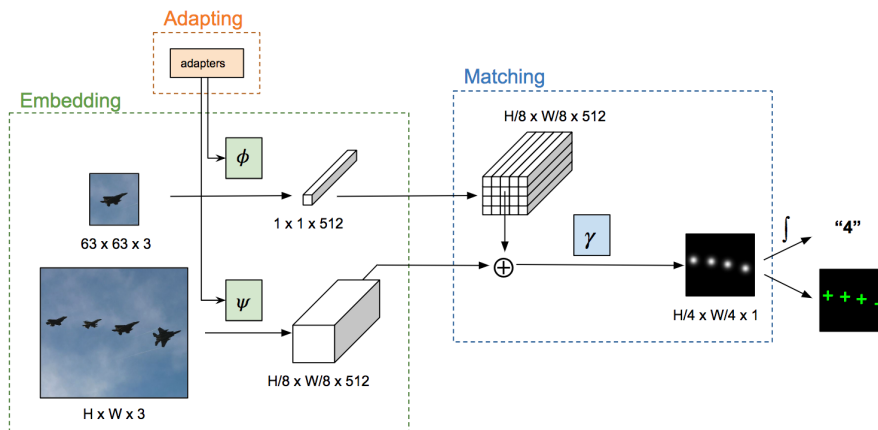


Fig. 2: The GMN architecture consists of three modules: embedding, matching, and adapting. The final count and detections are obtained by taking the integral and local maximums, respectively, over the output heatmap. The integral-based method is used for counting instances with significant overlap, while the local max is used where objects are well-separated and the positional information is of interest for further processing. The \oplus represents channel-wise concatenation.

Given a trained matching model, several factors can prevent it from generalizing perfectly onto the target domain: for instance, the image statistics can be very different from the training set (e.g. natural images vs. microscopy images), or the user requirements can be different (e.g. counting cars vs. counting only red cars). Efficient domain adaptation requires a module that can change the network activations with minimal effort (that is, minimal number of trainable parameters and very few training data). Thus, for the *adapting* stage, we incorporate residual adapter modules [28] to specialize the GMN to such needs. Adapting to the new counting task then merely involves freezing the majority of parameters in the generic matching network, and training the adapters (a small number of extra parameters) on the target domain with only a few labeled examples.

2.1 Embedding

In this module, a two-stream network is defined for transforming raw RGB images into high-level feature encodings. The two streams are parametrized by

separate functions for higher representation capacity:

$$v = \phi(z; \theta_1) \quad f = \psi(x; \theta_2)$$

In detail, the function ϕ transforms an exemplar image patch $z \in R^{63 \times 63 \times 3}$ to a feature vector $v \in R^{1 \times 1 \times 512}$, and ψ maps the full image $x \in R^{H \times W \times 3}$ to a feature map $f \in R^{H/8 \times W/8 \times 512}$. Both the vector v and feature maps f are L2 normalized along the feature dimensions. In practice, our choices for $\phi(\cdot; \theta_1)$ and $\psi(\cdot; \theta_2)$ are ResNet-50 networks [15] truncated after the final conv3_x layer. The resulting feature map from the image patch is globally max-pooled into the feature vector v .

2.2 Matching

The relations between the resulting feature vector and maps are modeled by a trainable function $\gamma(\cdot; \theta_3)$ that takes the concatenation of v and f as input, and outputs a similarity heat map, as shown in Figure 2. Before concatenation, v is broadcast to match the size of the feature maps to accommodate the fully convolutional feature, which allows for efficient modeling of the relations between the exemplar object and all other objects in the image. The similarity Sim is given by

$$Sim = \gamma([broadcast(v) : f]; \theta_3)$$

where “:” refers to concatenation, and $\gamma(\cdot; \theta_3)$ is parametrized by one 3×3 convolutional layer and one 3×3 convolutional transpose layer with stride 2 (for upsampling).

2.3 Training Generic Matching Networks

The generic matching network (consisting of embedding and matching modules) is trained on the ILSVRC video dataset. The ground truth label is a Gaussian placed at each instance location, multiplied by a scaling factor of 100, and a weighted MSE (Mean Squared Error) loss is used. Regressing a Gaussian allows the final count to be obtained by simply summing over the output similarity map, which in this sense doubles as a density map.

During training, the exemplar images are re-scaled to size 63×63 pixels, with the object of interest centered to fit the patch, and the larger 255×255 search image is taken as a crop centered around the scaled object (architecture details can be found in Table 1. More precisely, we always scale the search image according to the bounding box (w,h) of the exemplar objects, where the scale factor is obtained by solving $s \times h \times w = 63^2$. The input data is augmented with horizontal flips and small ($< 25^\circ$) rotations and zooms, and we sample both positive and negative pairs. In all subsequent experiments, the network has been pre-trained as described here.

Module	Exemplar Patch ($N \times 63 \times 63 \times 3$)	Image to Count ($N \times 255 \times 255 \times 3$)	Output Size
<i>Embedding</i>	conv, 7×7 , 64, stride 2	conv, 7×7 , 64, stride 2	$N \times 32 \times 32 \times 64$ $N \times 128 \times 128 \times 64$
	max pool, 3×3 , stride 2 $\begin{bmatrix} \text{conv}, 1 \times 1, 64 \\ \text{conv}, 3 \times 3, 64 \\ \text{conv}, 1 \times 1, 256 \end{bmatrix} \times 3$	max pool, 3×3 , stride 2 $\begin{bmatrix} \text{conv}, 1 \times 1, 64 \\ \text{conv}, 3 \times 3, 64 \\ \text{conv}, 1 \times 1, 256 \end{bmatrix} \times 3$	$N \times 16 \times 16 \times 256$ $N \times 64 \times 64 \times 256$
	$\begin{bmatrix} \text{conv}, 1 \times 1, 128 \\ \text{conv}, 3 \times 3, 128 \\ \text{conv}, 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} \text{conv}, 1 \times 1, 128 \\ \text{conv}, 3 \times 3, 128 \\ \text{conv}, 1 \times 1, 512 \end{bmatrix} \times 4$	$N \times 8 \times 8 \times 512$ $N \times 32 \times 32 \times 512$
	Global Maxpool	No Operation	$N \times 1 \times 1 \times 512$ $N \times 32 \times 32 \times 512$
<i>Matching</i>	Vector Broadcasting (32×32)	No Operation	$N \times 32 \times 32 \times 512$ $N \times 32 \times 32 \times 512$
	Feature Map Concatenation		$N \times 32 \times 32 \times 1024$
	Relation Module $\begin{bmatrix} \text{conv}, 3 \times 3, 256 \\ \text{convt}, 3 \times 3, 256 \end{bmatrix}$		$N \times 64 \times 64 \times 256$
	Prediction $\begin{bmatrix} \text{conv}, 3 \times 3, 1 \end{bmatrix}$		$N \times 64 \times 64 \times 1$

Table 1: Architecture of the generic matching networks. “convt” refers to convolutional transpose with stride 2.

Once trained on the tracking data, the model can be directly applied for detecting repetitions within an image. We show a number of example predictions in Figure 3. Note here, several interesting phenomena can be seen: *first*, as expected, the generic matching network has learned to match instances beyond a simplistic level; for instance, the animals are of different viewpoints, the bird in the fourth row is partially occluded, and the persons are not only partially occluded, but also in different shirts with substantial appearance variations; *second*, object overlaps can also be handled, as shown in the airplane cases; *third*, although the ImageNet training set is only composed of natural images, and none of the categories has a similar appearance or distribution to the HeLa cells, the generic matching network succeeds despite large appearance and shape variations which exist for cells. These results validate our idea of building a class-agnostic counting network. However, it is crucial to be able to easily *adapt* the pre-trained model to further specialize to new domains.

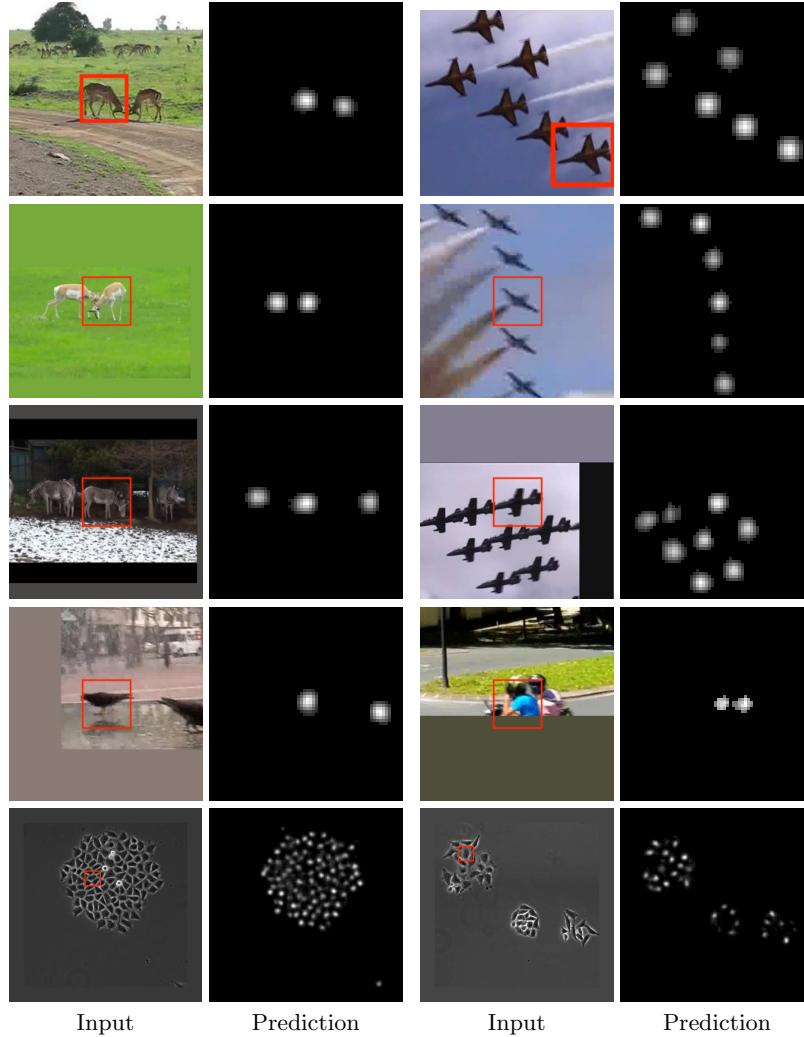


Fig. 3: Similarity predictions of the generic matching network on the video validation set, and on an *unseen* dataset of HeLa cells. The exemplar patch is marked with a red square. Images are padded with the mean value and the resolution has been changed for visualization purposes. As expected, the generic matching network has learned to match instances beyond a simplistic level; for instance, the animals are of different viewpoints, the matched bird in the fourth row is partially occluded, and the people are in different colored shirts. More interestingly, it acts as an excellent initialization for objects from unseen domains, even in the presence of large appearance and shape variation in the case of HeLa cells.

2.4 Adapting

The next objective is to specialize the network to new domains or new user requirements. We add *residual adapter modules* [28] implemented as 1×1 convolutions in parallel with the existing 3×3 convolutions in the embedding module of the network. During adaptation, we freeze all of the parameters in the pre-trained generic matching network, and train only the adapters and batch normalization layers. This results in 178K trainable parameters out of a total network size of 6.0M parameters, only 3% of the total.

2.5 Discussion and relation to prior work

Object counting poses certain additional challenges that are less prominent or non-existent in tracking. First, rather than requiring a single maximum in a candidate window (that localizes the object), counting requires a clean output map to distinguish multiple matches from noise and false positives. Second, unlike the continuous variation of object shape and appearance in the tracking problem, object counting can have more challenging appearance changes, e.g. large degrees of rotation, and intra-class variation (in the case of cars, both color and shape). Thus, we find the approaches used in template matching (SSD or cross-correlation [6,9,22]) to be insufficient for our purposes (as will be shown in Table 4). To address these challenges, we learn a discriminative classifier $\gamma(\cdot; \theta_3)$ between the exemplar patch and search image, an idea that dates back to [23].

The residual adapters [28] are added only to the embedding module, but we train the batch normalization layers throughout the entire network. Marsden et al. [25] also use residual adapters to adapt a network for counting different objects. However, they place the modules in the final fully connected layers in order to regress a count, whereas we add them to the convolutional layers in the residual blocks of the embedding module, such that they are able to change the filter responses at every stage of the base model, providing more capacity for adaptation.

3 Counting Benchmark Experiments

As a proof of concept, the generic matching network is visually validated as a strong initialization for counting objects from unseen domains (Figure 3). To further demonstrate the effectiveness of the general-purpose GMN, we adapt the network to three different datasets: VGG synthetic cells [20,21], HeLa phase-contrast cells [3], and a large-scale drone-collected car dataset [16].

Each of these datasets poses unique challenges. The synthetic cells contain many overlapping instances, a condition where density estimation methods have shown strong performance. The HeLa cells exhibit significantly more variation in size and appearance than the synthetic cells, and the number of training images is extremely limited (only 11 images); thus, detection-based methods with handcrafted features have shown good results. In the car dataset, cars appear in

various orientations, often within the same image, and can be partially occluded by trees and bridges; there is also clutter from motorbikes, buildings, and other distractors (Figure 6). As shown in Hsieh et al. [16], state-of-the-art models for object detection produce a very high error rate.

3.1 Evaluation Metrics

The metrics we use for evaluation throughout this paper are the mean absolute counting error (MAE), precision, recall, and F_1 score. To determine successful detections, we first take the local maximums (above a threshold T) of the predicted similarity map as the detections. T is usually set as the value that maximizes the F_1 score on a validation set. Note that, since multiple combinations of recall and precision can give the same F_1 score, we prioritize the recall score. Following [3], we then match these predicted detections with the ground truth locations using the Hungarian algorithm, with the constraint that a successful detection must lie no further than a tolerance R from the ground truth location, where R is set as the average radius of each object.

3.2 Synthetic fluorescence microscopy

The synthetic VGG cell dataset contains 200 fluorescence microscopy cell images, evenly split between training and testing sets. We follow the procedure proposed by Lempitsky and Zisserman [21] of sampling 5 random splits of the training set with N training images and N validation images. Results in Table 2 and Figure 4 show that our method is not restricted to detection-based counting, but also performs well on density estimation-type problems in a setting with high instance overlap. Note that, we compare with methods that are highly engineered for this dataset.

Method	MAE	Precision	Recall	F_1 -score
Xie et al.[34]	2.9 ± 0.2	-	-	-
Fiaschi et al.[12]	3.2 ± 0.1	-	-	-
Lempitsky and Zisserman [21]	3.5 ± 0.2	-	-	-
Barinova et al.[5]	6.0 ± 0.5	-	-	-
Singletons [3]	51.2 ± 0.8	98.87 ± 1.52	72.07 ± 0.85	83.37 ± 1.20
Full system w/o surface [3]	5.06 ± 0.2	95.00 ± 0.75	91.97 ± 0.43	93.46 ± 0.15
Ours	3.56 ± 0.27	99.43 ± 0.05	82.50 ± 0.15	90.18 ± 0.07

Table 2: Results for the synthetic cell dataset. All methods are trained on the $N = 32$ split. Standard deviations are calculated using 5 random splits of training and validation sets and 5 randomly sampled exemplar patches per image. Note here, the exemplar patches are sampled from images in the training set, and different exemplar patches have negligible effect on performance.

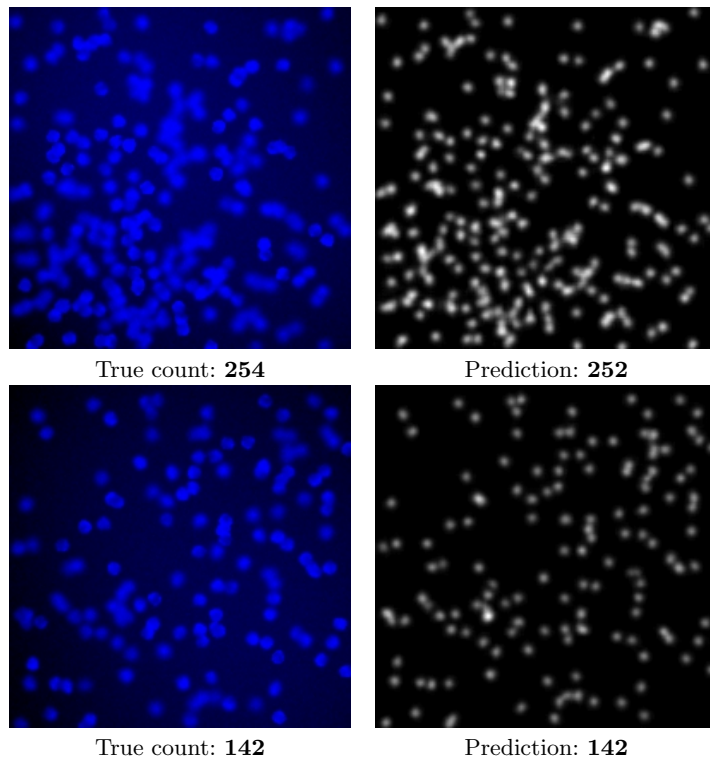


Fig. 4: Example of counting results on synthetic cell images. For each pair of images, left: original image, and right: the network’s predicted heat map, which is *summed* to give the estimated count.

3.3 HeLa cells on phase contrast microscopy

The dataset contains 11 training and 11 testing images. We follow the training procedure of [3] and train in a leave-one-out fashion for selecting hyperparameters, e.g. detection threshold T . Results are shown in Table 3. As shown in Figure 5, our method performs well in scenarios of large intra-class variations in shape and size, where SSD and cross-correlation would suffer. Overall, our GMN achieves comparable results to the conventional methods with hand-crafted features, despite the training dataset being extremely small for current deep learning standards.

3.4 Cars

We next demonstrate the GMN’s performance on counting cars in aerial images. This drone-collected dataset (CARPK) consists of 989 training images and 459

Method	MAE	Precision	Recall	F1-score
Correlation clustering [36]	-	-	-	95
Singletons [3]	2.36 ± 0.67	93.70 ± 0.20	91.94 ± 0.72	92.81 ± 0.35
Full system w/o surface [3]	3.84 ± 1.44	98.51 ± 1.16	95.76 ± 0.27	97.10 ± 0.27
Ours	3.53 ± 0.18	96.05 ± 0.04	94.22 ± 0.06	95.12 ± 0.05

Table 3: Results for the HeLa cell dataset. We calculate MAE using the detection counts, since the instances are well-separated. Our standard deviations are calculated using 5 randomly sampled exemplar patches per image. Note here, the 5 exemplar patches are sampled from images in the training set; different exemplar patches have negligible effect on performance.

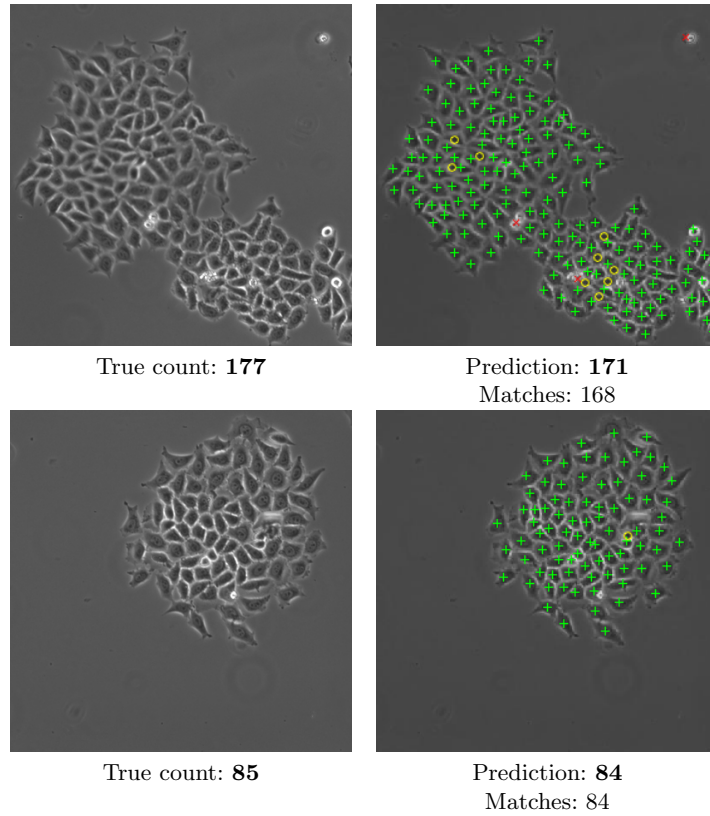


Fig. 5: Example detection results on the HeLa cell test set. Correct detections (based on Hungarian matching) are marked with a green '+', false positives with a red 'x', and missed detections with a yellow 'o'.

testing images (nearly 90,000 instances of cars), where the images are taken from overhead shots of car parking lots. The training images are taken from three different parking lot scenes, and the test set is taken from a fourth scene. We compare our network to the region proposal and classification methods in Table 4.

In the experiments, we train two GMN models with augmentation: one on just three images (99 total cars) randomly sampled from the training scenes, which achieves state-of-the-art results, and one on the full CARPK training set, which further boosts the performance by a large margin.

Method	T	MAE	RMSE	Recall	Precision
*YOLO [16,29]	-	48.89	57.55	-	-
*Faster R-CNN [16,30]	-	47.45	57.39	-	-
*Faster R-CNN (RPN-small) [16,30]	-	24.32	37.62	-	-
†One-Look Regression [16,26]	-	59.46	66.84	-	-
*Spatially Regularized RPN [16]	-	23.80	36.79	57.5%	-
Template matching (Sum of Squared Distances)	-	49.8	59.7	20.0%	29.1%
Ours (3 images, 99 cars)	2.5	36.71	44.16	60.65%	93.91%
Ours (3 images, 99 cars)	2	22.32	28.72	71.32%	90.23%
Ours (3 images, 99 cars)	1.75	17.32	22.81	74.16%	87.87%
Ours (3 images, 99 cars)	1.5	13.38	18.03	76.1%	85.1%
Ours (full dataset)	2.75	19.66	25.12	78.61%	97.0%
Ours (full dataset)	2.5	14.36	19.01	83.2%	96.0%
Ours (full dataset)	2	8.38	11.55	87.46%	93.4%
Ours (full dataset)	1.75	7.48	9.9	88.4%	91.8%

Table 4: Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), Recall and Precision comparisons on the CARPK dataset. The “*” indicates that the method has been fine tuned on the full dataset, and the “†” indicates that the method has been revised to fit the dataset, as described in [16]. We show our method trained on 3 images and on the full dataset, with varying thresholds T . We calculate MAE using 5 randomly sampled exemplar patches per image, and the final counts are obtained from local maximums (counting by detection). Note here, the exemplar patches are sampled from images in the training set, and different exemplar patches have negligible effect on performance. Standard deviation is not reported in this table, but can be easily computed following the previously reported manner.

When determining counts based on local maximums, we note it is possible that our model outperforms the previous detection-based methods due to false positives and false negatives “canceling” each other, making the counting error very low. Thus, we investigate effects of the threshold T (as defined in § 3.1) on selecting detections from candidate local maximums, and report results for several values of T . Note that, by varying this hyperparameter, we are able to

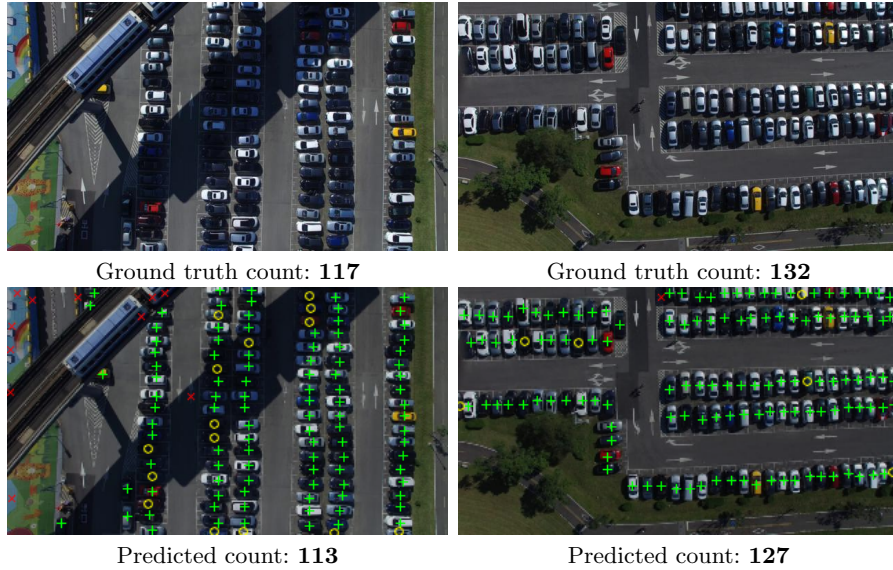


Fig. 6: Sample results on the CARPK dataset. Top row: original images. Bottom row: predicted detections. Correct detections are marked with a green ‘+’, false positives with a red ‘x’, and missed detections with a yellow ‘o’. Many of the missed detections are dark cars in shadow, which upon inspection are difficult for even a human eye to discern.

explicitly control the precision-recall of our model. While calculating recall and precision, we consider a detection to be successful if it lies within 20 pixels (determined based on the mean car size) of the ground truth location. The recall reported for the region proposal methods in Table 4 is calculated by averaging across scores from using various IoU thresholds, as described in [16].

As shown in Table 4, the MAE is calculated with 5 randomly sampled exemplar patches per image, and the final counts are obtained by counting local maximums (detection-based counting). Note here, the exemplar patches are sampled from images in the training set, and different exemplar patches have negligible effect on performance. We can see that even with a very high precision (model trained on the full dataset, with $T = 2.75$), our model can still outperform the previous state-of-the-art by a substantial margin (counting error: MAE=23.8 vs MAE=19.7). Further decreasing the threshold yields higher recall at the expense of precision, with our best model achieving a counting error of MAE=7.5.

3.5 Discussion

From our experiments, the following phenomena can be observed:

First, in contrast to previous work, where different architectures are designed for density estimation in scenarios with significant instance overlap (e.g. VGG synthetic cells) and for detection-based counting in scenarios with well-separated objects (e.g. HeLa cells and cars), the GMN has the flexibility to handle both scenarios. Based on the amount of instance overlap, the object counts can simply be obtained by taking either the integral in the former case, or the local maximum in the latter, or possibly even an ensemble of them [17].

Second, by training in a discriminative manner, the GMN is able to match instances beyond the simplistic level, making it more robust to large degrees of rotation and appearance variation than the baseline SSD-based template matching (as shown in Table 4).

Third, in the cases where training data is limited (11 images for HeLa cells, 3 for cars), the proposed model has consistently shown comparable or superior performance to the state-of-the-art methods, indicating the model’s ease of adaptation, as well as verifying our observation that videos can be a natural data source for learning *self-similarity*.

4 ShanghaiTech Crowd Counting

To further demonstrate the power and flexibility of counting-by-matching, we extend it to the ShanghaiTech crowd counting dataset, which contains images of very large crowds of people from arbitrary camera perspectives, with individuals appearing at extremely varied scales due to perspective.

We carry out a preliminary implementation of our method on the ShanghaiTech Part A crowd dataset. Inspired by the idea of crowd detection as repetitive textures [1], we conjecture that it is possible to ignore individual instances and match the statistics of patches instead; e.g. the statistics of patches with 10 people should be different from those with 20 people.

We take the following steps: (1) Using the ground truth dot annotations, we quantize 64×64 pixel patches into 10 different classes based on number of people, e.g. one class will be 0 people, another 5 people, etc. (See Figure 7 for an example.) (2) Following the idea of counting-by-matching, we train the self-similarity architecture to embed the patches based on the number of people, i.e. if patches are sampled from the same class, the model must predict 1, otherwise 0. (3) We run the model on the test set using a sample of each class from the training set as the exemplar patch, with the final classification made by the maximum response.

Compared to other models that are specifically designed to count human crowds (e.g. CNNs with multiple branches), we aim for a method with the potential for low-shot category-agnostic counting. Our preliminary experiments show the possibility of scaling the counting-by-matching idea to human crowd datasets.

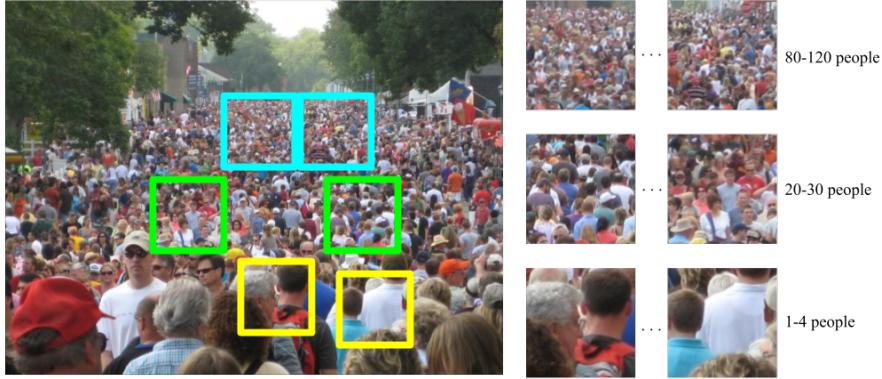


Fig. 7: Example of how different patches are classified. Patches marked by the same color square belong to the same density class. As can be seen, the significant textural differences between the classes enables a network to learn to classify different densities.

Method	MAE	RMSE
†Zhang et al. [35]	181.8	277.7
MCNN-CCR [37]	245.0	336.1
MCNN [37]	110.2	173.2
ic-CNN [27]	69.8	117.3
Ours	95.8	133.3

Table 5: Preliminary results on Shanghaitech Part A (lower is better). Patches chosen based on validation set performance. The “†” result is from paper [37].

5 Conclusion

In this work, we recast counting as a matching problem, which offers several advantages over traditional counting methods. Namely, we make use of object detection video data that has not yet been utilized by the counting community, and we create a model that can flexibly adapt to various domains, which is a form of few-shot learning. We hope this unconventional structuring of the counting problem encourages further work towards an all-purpose counting model.

Several extensions are possible for future works: *first*, it would be interesting to consider counting in video sequences, rather than individual images or frames. Here the tracking analogue takes on an even greater significance as a counting model can take advantage of both within-frame and between-frame similarities, *second*, a carefully engineered scale-invariant network with more sophisticated feature fusion than the GMN could potentially improve the current results.

Acknowledgements

Funding for this research is provided by the Oxford-Google DeepMind Graduate Scholarship, and by the EPSRC Programme Grant Seebibyte EP/M013774/1.

References

1. Arandjelovic, O.: Crowd detection from still images. In: Proc. BMVC. (2008)
2. Arteta, C., Lempitsky, V., Noble, J.A., Zisserman, A.: Interactive object counting. In: Proc. ECCV (2014)
3. Arteta, C., Lempitsky, V., Noble, J.A., Zisserman, A.: Detecting overlapping instances in microscopy images using extremal region trees. *Medical Image Analysis* (2015)
4. Arteta, C., Lempitsky, V., Zisserman, A.: Counting in the wild. In: Proc. ECCV (2016)
5. Barinova, O., Lempitsky, V., Kohli, P.: On the detection of multiple object instances using Hough transforms. In: Proc. CVPR (2010)
6. Bertinetto, L., Valmadre, J., Henriques, J.F., Vedaldi, A., Torr, P.H.S.: Fully-convolutional siamese networks for object tracking. In: Workshop on Visual Object Tracking, ECCV (2016)
7. Buades, A., Coll, B., Morel, J.M.: A non-local algorithm for image denoising. In: Proc. CVPR. pp. 60–65 (2005)
8. Cho, S., Chow, T., Leung, C.: A neural-based crowd estimation by hybrid global learning algorithm. In: IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics) (2009)
9. Dekel, T., Oron, S., Rubinstein, M., Avidan, S., Freeman, W.: Best-buddies similarity for robust template matching. In: Proc. CVPR (2015)
10. Desai, C., Ramanan, D., Fowlkes, C.: Discriminative models for multi-class object layout. In: Proc. ICCV (2009)
11. Efros, A., Leung, T.: Texture synthesis by non-parametric sampling. In: Proc. ICCV. pp. 1039–1046 (Sep 1999)
12. Fiaschi, L., Nair, R., Köethe, U., Hamprecht, F.: Learning to count with regression forest and structured labels. In: Proc. ICPR (2012)
13. Glasner, D., Bagon, S., Irani, M.: Super-resolution from a single image. In: Proc. ICCV (2009)
14. Han, X., Leung, T., Jia, Y., Sukthankar, R., Berg, A.: Matchnet: Unifying feature and metric learning for patch-based matching. In: Proc. CVPR (2015)
15. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proc. CVPR (2016)
16. Hsieh, M., Lin, Y., Hsu, W.: Drone-based object counting by spatially regularized regional proposal networks. In: Proc. ICCV (2017)
17. Idrees, H., Tayyab, M., Athrey, K., Zhang, D., Al-Máadeed, S., Rajpoot, N., Shah, M.: Composition loss for counting, density map estimation and localization in dense crowds. In: Proc. ECCV (2018)
18. Koch, G., Zemel, R., Salakhutdinov, R.: Siamese neural networks for one-shot image recognition. In: ICML 2015 Deep Learning Workshop (2015) (2015)
19. Kong, D., Gray, D., Tao, H.: A viewpoint invariant approach for crowd counting. In: Proc. ICPR. vol. 3, pp. 1187–1190. IEEE (2006)

20. Lehmussola, A., Ruusuvuori, P., Selinummi, J., Huttunen, H., Yli-Harja, O.: Computational framework for simulating fluorescence microscope images with cell populations. *IEEE Transactions on Medical Imaging* (2007)
21. Lempitsky, V., Zisserman, A.: Learning to count objects in images. In: *NIPS* (2010)
22. Leung, T., Malik, J.: Detecting, localizing and grouping repeated scene elements from an image. In: *Proc. ECCV* (1996)
23. Malisiewicz, T., Gupta, A., Efros, A.A.: Ensemble of exemplar-SVMs for object detection and beyond. In: *Proc. ICCV* (2011)
24. Marana, A., Velastin, S., Costa, L., Lotufo, R.: Estimation of crowd density using image processing. In: *Image Processing for Security Applications*. pp. 11–1 (1997)
25. Marsden, M., McGuinness, K., S., L., Keogh, C.E., O’Connor, N.E.: People, penguins and petri dishes: Adapting object counting models to new visual domains and object types without forgetting. In: *Proc. CVPR* (2018)
26. Mundhenk, T.N., Konjevod, G., Sakla, W.A., Boakye, K.: A large contextual dataset for classification, detection and counting of cars with deep learning. In: *Proc. ECCV* (2014)
27. Ranjan, V., Le, H., Hoai, M.: Iterative crowd counting. In: *Proc. ECCV* (2018)
28. Rebuffi, S.A., Bilen, H., Vedaldi, A.: Efficient parametrization of multi-domain deep neural networks. In: *Proc. CVPR* (2018)
29. Redmon, J., Divvala, S.K., Girshick, R.B., Farhadi, A.: You only look once: Unified, real-time object detection. In: *Proc. CVPR* (2016)
30. Ren, S., He, K., Girshick, R., Sun, J.: Faster R-CNN: Towards real-time object detection with region proposal networks. In: *NIPS* (2016)
31. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, S., Karpathy, A., Khosla, A., Bernstein, M., Berg, A., Li, F.: Imagenet large scale visual recognition challenge. *IJCV* (2015)
32. Sung, F., Yang, Y., Zhang, L., Xiang, T., Torr, P.H.S., Hospedales, T.: Learning to compare: Relation network for few-shot learning. In: *Proc. CVPR* (2018)
33. Vinyals, O., Blundell, C., Lillicrap, T., Kavukcuoglu, K., Wierstra, D.: Matching networks for one shot learning. In: *NIPS* (2016)
34. Xie, W., Noble, J.A., Zisserman, A.: Microscopy cell counting with fully convolutional regression networks. In: *MICCAI 1st Workshop on Deep Learning in Medical Image Analysis* (2015)
35. Zhang, C., Li, X., Wang, X., Yang, X.: Cross-scene crowd counting via deep convolutional neural networks. In: *Proc. CVPR* (2015)
36. Zhang, C., Yarkony, J., Hamprecht, F.A.: Cell detection and segmentation using correlation clustering. In: *MICCAI* (2014)
37. Zhang, Y., Zhou, D., Chen, S., Gao, S., Ma, Y.: Single-image crowd counting via multi-column convolutional neural network. In: *Proc. CVPR* (2016)

8

Conclusion

Here we summarize the main contributions of this thesis and its impact thus far (Section 8.1). We then discuss ongoing and potential extensions to this work in Section 8.2.

8.1 Achievements and Impact

Grouping people with their correlated effects. In Chapter 3, we introduced a novel method for the new task of decomposing a video into layers containing one or more people and their correlated visual effects. The method overfits a neural network to a single video and requires no human annotation, instead taking advantage of the deep image prior property of CNNs to separate the scene elements. We additionally train a CNN that converts human keypoints to UV maps using pseudo-ground truth labels generated from an existing dataset and our own data (our ‘keypoint-to-UV network’). Our method first preprocesses the video to automatically detect and track the keypoints of all people and converts them to UV maps using the pretrained keypoint-to-UV network. We demonstrated the quality of the layer decomposition by using the layers to produce a variety of editing effects: removing individual people, retiming them (slowing / speeding up motions, and aligning the jumps of individuals), and duplicating people. We tested our method on videos containing

complex scene elements, such as water splashes, trampoline deformations, and group running. The code was open-sourced on Github and has been starred over 100 times in the 8 months since it has been released.¹

Grouping objects with their correlated effects. In Chapter 4, we built upon our work in Chapter 3 to present a method that can produce layers for *any* object class, not solely humans. We accomplished this generalization by showing that UV geometry is not a necessary input to the network, and that it can be replaced by rough binary object masks and precomputed optical flow to provide correspondence in lieu of geometry. We additionally introduced flow-based losses to provide greater temporal coherence and improved layer separation, which we showed in ablation studies. We demonstrated results on a wide range of object classes, including animals, vehicles, and sports equipment. We additionally generated editing effects using our predicted layers, including stroboscopy, color manipulation, and background replacement. This work was selected for oral presentation at the *Conference on Computer Vision and Pattern Recognition*, and the code has been released on Github.² The code has been starred over 450 times in the 2 months since release, demonstrating its value to the research community.

Self-supervised video object segmentation. In Chapter 5, we presented a method for densely tracking object segments in video, and trained the model in a fully self-supervised manner, i.e. without requiring any human annotation. The model incorporates an efficient memory module which enables it to overcome difficult occlusions of the tracked object. Our method significantly outperformed the state-of-the-art self-supervised methods on multiple standard benchmarks, narrowing the gap with supervised methods, and displaying better generalization ability to unseen object classes than supervised methods. The code and trained model have been released on Github, and the repository³ has been starred 247

¹The code is available at: <https://github.com/google/retiming/>.

²The code is available at: <https://github.com/erikalau/omnimatte/>.

³The code and trained model is available at: <https://github.com/zlai0/MAST>.

times in the 14 months since release. The paper has garnered 42 citations in the same amount of time. Several subsequent works have since improved on our state-of-the-art results [A. A. Jabri et al. 2020; J. Xu and X. Wang 2021], indicating interest in this research area.

Self-supervised moving object discovery. In Chapter 6, we presented a self-supervised method for discovering and segmenting moving objects in video. The model operates solely on precomputed optical flow inputs; it decomposes the input into two optical flow layers, one representing the background motion and the other representing the foreground object motion. The method performs at par or better than previous self-supervised methods on existing benchmarks, while running an order of magnitude faster. We additionally demonstrate our model’s ability to segment challenging camouflaged objects by exploiting motion cues rather than misleading RGB information. The code has been open-sourced⁴ and has been starred over 100 times within the first 3 months of its release. An abbreviated version of the paper received the Best Paper Award at the *Workshop on Robust Video Scene Understanding: Tracking and Video Segmentation* at the *Conference on Computer Vision and Pattern Recognition*.

Class-agnostic counting. In Chapter 7, we take a brief detour from video understanding and instead use videos as a source of training data for an unlikely task: counting. While existing methods are designed to count a specific object class, we observe that humans are able to easily count objects of never-before seen object classes. Thus we propose the first class-agnostic counting method, which we pretrain on video data labeled for tracking. The model is able to quickly adapt to new domains (such as cell biology) using small amounts of labeled data, despite having been pretrained on visually distinct data (Imagenet Video [Russakovsky et al. 2015]). This work has been cited 38 times and has inspired follow-up work such as SIMCO

⁴The code is available at: <https://github.com/charigyang/motiongrouping>.

[Godi et al. 2021], and Dwibedi et al. 2020, which applies the same principle towards counting repeated actions in video. The code has been open-sourced on Github⁵.

8.2 Future Work

8.2.1 Faster Omnimatte

The Omnimatte method optimizes the parameters of a neural network to overfit to a single video. The method requires several hours to produce a result for a single video—roughly 2 hours for a video with 1 object and 82 frames, and this number increases with the length of the video and number of objects.

Data-Driven Omnimatte. One avenue to improving the performance is to train a single network on a corpus of videos to perform the layer decomposition in a single feed-forward pass. To apply a pretrained model to a new video requires significant restructuring of the network and inputs: rather than inputting a single layer’s object mask and outputting the RGBA for that layer, the network will need to observe the *original video frames*. Thus the layer separation cannot rely on the deep image prior property, as in the original Omnimatte work. The network will require some form of supervision or information bottleneck to encourage a meaningful decomposition of the RGB frames into individual object layers. One option is to run the original Omnimatte method on a large corpus of videos, and use the layer results as pseudo-ground truth to train a supervised generalized Omnimatte. This data may be supplemented by synthetic training data. Large video object segmentation datasets such as YouTubeVOS [N. Xu et al. 2018] and YouTubeVIS [Yang et al. 2019] contain manually annotated object segments which can be used to generate synthetic layers. Synthetic effects such as shadows can be simulated and composited along with the object segment, and can be designed to be correlated in shape and movement with the original object.

⁵The code is available at: <https://github.com/erikaluc/class-agnostic-counting>.

Omnimatte with learned gradient descent. Another avenue for improving the performance of the Omnimatte method is a hybrid approach between the original optimization-based Omnimatte method and the previously proposed data-driven, single feed-forward pass method. We can aim for a small number of test-time optimization steps by training a model to predict several gradient updates to the Omnimatte network, similarly to the method used in Flynn et al. 2019 to speed up the original work on multi-plane images (depth layers) [T. Zhou et al. 2018]. Flynn et al. 2019 replace the gradient descent updates with a deep network that predicts the parameter updates. This deep network is trained on many examples, thus allowing for a learned prior on the model parameters. Their method results in convergence in just a few iterations, making it much faster than the original work by T. Zhou et al. 2018. We can apply this framework to the Omnimatte method to obtain similar improvements in speed.

8.2.2 End-to-End Training

This thesis considered several fundamental problems in computer vision, including object discovery, object segmentation, and layer decomposition for videos. One interesting direction is to unify these various self-supervised components into a single system that can be trained self-supervised end-to-end. Furthermore, while depth does not fall under the scope of this thesis, incorporating depth estimation and prediction in future work can further improve results. A unified system which performs a range of visual processing tasks such as depth estimation, object segmentation, and layer decomposition, and is trained in a self-supervised manner, would provide an elegant solution to fundamental computer vision problems.

Bibliography

- Edward H. Adelson (1991). *Layered representation for image coding*. Tech. rep. 181. The MIT Media Lab.
- Jonas Adler and Ozan Oktem (2017). “Solving ill-posed inverse problems using iterative deep neural networks”. In: *Inverse Problems*.
- (2018). “Learned primal-dual reconstruction.” In: *IEEE Transactions on Medical Imaging*.
- Triantafyllos Afouras, Yuki M. Asano, Francois Fagan, Andrea Vedaldi, and Florian Metze (2021). *Self-supervised object detection from audio-visual correspondence*. arXiv: [2104.06401](https://arxiv.org/abs/2104.06401) [cs.CV].
- Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman (2018a). “Deep Lip Reading: a comparison of models and an online application”. In: *INTERSPEECH*.
- (2018b). “The Conversation: Deep Audio-Visual Speech Enhancement”. In: *INTERSPEECH*.
- (2020). “ASR is all you need: Cross-modal distillation for lip reading”. In: *International Conference on Acoustics, Speech, and Signal Processing*.
- Pulkit Agrawal, Joao Carreira, and Jitendra Malik (2015). “Learning to see by moving”. In: *Proceedings of the International Conference on Computer Vision*. IEEE, pp. 37–45.
- Jean-Baptiste Alayrac, João Carreira, Relja Arandjelovic, and Andrew Zisserman (2019). “Controllable Attention for Structured Layered Video Decomposition”. In: *Proceedings of the International Conference on Computer Vision*.
- Jean-Baptiste Alayrac, João Carreira, and Andrew Zisserman (2019). “The Visual Centrifuge: Model-Free Layered Video Representations”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky (2020). “Neural Point-Based Graphics”. In: *Proceedings of the European Conference on Computer Vision*.
- Humam Alwassel, Dhruv Mahajan, Lorenzo Torresani, Bernard Ghanem, and Du Tran (2019). “Self-Supervised Learning by Cross-Modal Audio-Video Clustering”. In: *arXiv preprint arXiv:1911.12667*.

- Relja Arandjelović and Andrew Zisserman (2017). “Look, Listen and Learn”. In: *Proceedings of the International Conference on Computer Vision*.
- (2018). “Objects that sound”. In: *Proceedings of the European Conference on Computer Vision*.
- Simon Baker, Richard Szeliski, and P. Anandan (1998). “A Layered Approach to Stereo Reconstruction”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- David Bau, Jun-Yan Zhu, Hendrik Strobelt, Bolei Zhou, Joshua B. Tenenbaum, William T. Freeman, and Antonio Torralba (2019). “GAN Dissection: Visualizing and Understanding Generative Adversarial Networks”. In: *Proceedings of the International Conference on Learning Representations*.
- Sagie Benaim, Ariel Ephrat, Oran Lang, Inbar Mosseri, William T. Freeman, Michael Rubinstein, Michal Irani, and Tali Dekel (2020). “SpeedNet: Learning the Speediness in Videos”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Alexander Bokov and Dmitriy Vatolin (2018). “100+ Times Faster Video Completion by Optical-Flow-Guided Variational Refinement”. In: *IEEE International Conference on Image Processing (ICIP)*.
- Gabriel J Brostow and Irfan A Essa (1999). “Motion based decompositing of video”. In: *Proceedings of the International Conference on Computer Vision*.
- Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros (2019). “Everybody dance now”. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5933–5942.
- Ya-Liang Chang, Zhe Yu Liu, Kuan-Ying Lee, and Winston Hsu (2019). “Free-form Video Inpainting with 3D Gated Convolution and Temporal PatchGAN”. In: *Proceedings of the International Conference on Computer Vision*.
- Changan Chen, Unnat Jain, Carl Schissler, Sebastia Vicenc Amengual Gari, Ziad Al-Halah, Vamsi Krishna Ithapu, Philip Robinson, and Kristen Grauman (2020). “Soundspaces: Audio-visual navigation in 3d environments”. In: *Proceedings of the European Conference on Computer Vision*.
- Qifeng Chen and Vladlen Koltun (2017). “Photographic image synthesis with cascaded refinement networks”. In: *Proceedings of the International Conference on Computer Vision*.
- Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman (2017). “Lip Reading Sentences in the Wild”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Joon Son Chung and Andrew Zisserman (2016). “Lip Reading in the Wild”. In: *Proceedings of the Asian Conference on Computer Vision*.

- Trevor Darrell and Alex Pentland (1991). “Robust estimation of a multi-layered motion representation.” In: *IEEE Workshop on Visual Motion*.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei (2009). “ImageNet: A Large-Scale Hierarchical Image Database”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Emily Denton and Vighnesh Birodkar (2017). “Unsupervised learning of disentangled representations from video”. In: *Advances in Neural Information Processing Systems*.
- Ali Diba, Vivek Sharma, Luc Van Gool, and Rainer Stiefelhagen (2019). “DynamoNet: Dynamic Action and Motion Network”. In: *Proceedings of the International Conference on Computer Vision*.
- Debidatta Dwibedi, Yusuf Aytar, Jonathan Tompson, Pierre Sermanet, and Andrew Zisserman (2020). “Counting Out Time: Class Agnostic Video Repetition Counting in the Wild”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Dave Epstein, Boyuan Chen, and Carl Vondrick (2020). “Oops! Predicting Unintentional Action in Video”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Basura Fernando, Hakan Bilen, Efstratios Gavves, and Stephen Gould (2017). “Self-Supervised Video Representation Learning With Odd-One-Out Networks”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- John Flynn, Michael Broxton, Paul Debevec, Matthew DuVall, Graham Fyffe, Ryan Overbeck, Noah Snavely, and Richard Tucker (2019). “DeepView: View synthesis with learned gradient descent”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Matthieu Fradet, Patrick Pérez, and Philippe Robert (2008). “Semi-automatic Motion Segmentation with Motion Layer Mosaics”. In: *Proceedings of the European Conference on Computer Vision*.
- Chuang Gan, Boqing Gong, Kun Liu, Hao Su, and Leonidas J. Guibas (2018). “Geometry guided convolutional neural networks for self-supervised video representation learning”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Yossi Gandelsman, Assaf Shocher, and Michal Irani (2019). ““Double-DIP”: Unsupervised Image Decomposition via Coupled Deep-Image-Priors”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Chen Gao, Ayush Saraf, Jia-Bin Huang, and Johannes Kopf (2020). “Flow-edge Guided Video Completion”. In: *Proceedings of the European Conference on Computer Vision*.
- Marco Godi, Christian Joppi, Andrea Giachetti, and Marco Cristani (2021). “SIMCO: SIMilarity-based object COunting”. In: *Proceedings of the International Conference on Pattern Recognition*.

- Tengda Han, Weidi Xie, and Andrew Zisserman (2019). “Video Representation Learning by Dense Predictive Coding”. In: *Workshop on Large Scale Holistic Video Understanding, ICCV*.
- (2020a). “Memory-augmented Dense Predictive Coding for Video Representation Learning”. In: *Proceedings of the European Conference on Computer Vision*.
- (2020b). “Self-supervised Co-training for Video Representation Learning”. In: *Advances in Neural Information Processing Systems*.
- David Harwath, Adria Recasens, Didac Suris, Galen Chuang, Antonio Torralba, and James Glass (2018). “Jointly Discovering Visual Objects and Spoken Words from Raw Sensory Input”. In: *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow (2018). “Deep Blending for Free-viewpoint Image-based Rendering”. In: *ACM Trans. Graph.* 37.6, 257:1–257:15.
- Jia-Bin Huang, Sing Bing Kang, Narendra Ahuja, and Johannes Kopf (2016). “Temporally coherent completion of dynamic video”. In: *tog*.
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros (2017). “Image-to-Image Translation with Conditional Adversarial Networks”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Phillip Isola, Daniel Zoran, Dilip Krishnan, and Edward H Adelson (2015). “Learning visual groups from co-occurrences in space and time”. In: *Proceedings of the International Conference on Learning Representations*.
- Allan A. Jabri, Andrew Owens, and Alexei A. Efros (2020). “Space-Time Correspondence as a Contrastive Random Walk”. In: *Advances in Neural Information Processing Systems*.
- Tomas Jakab, Ankush Gupta, Hakan Bilen, and Andrea Vedaldi (2018). “Conditional image generation for learning the structure of visual objects”. In: *Advances in Neural Information Processing Systems*.
- (2020). “Self-supervised Learning of Interpretable Keypoints from Unlabelled Videos”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Dinesh Jayaraman and Kristen Grauman (2015). “Learning image representations tied to ego-motion”. In: *Proceedings of the International Conference on Computer Vision*.
- (2016). “Slow and steady feature analysis: higher order temporal coherence in video”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Longlong Jing, Xiaodong Yang, Jingen Liu, and Yingli Tian (2018). “Self-supervised Spatiotemporal Feature Learning by Video Geometric Transformations”. In: *arXiv preprint arXiv:1811.11387*.

- Njegica Jojic and B.J. Frey (2001). “Learning flexible sprites in video layers”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman (2017). “The Kinetics Human Action Video Dataset.” In: *arXiv preprint arXiv:1705.06950*.
- Dahun Kim, Donghyeon Cho, and In So Kweon (2018). “Self-Supervised Video Representation Learning with Space-Time Cubic Puzzles”. In: *AAAI*.
- Dahun Kim, Sanghyun Woo, Joon-Young Lee, and In So Kweon (2019). “Deep Video Inpainting”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Bruno Korbar, Du Tran, and Lorenzo Torresani (2018). “Cooperative Learning of Audio and Video Models from Self-Supervised Synchronization”. In: *Advances in Neural Information Processing Systems*.
- Hilde Kuehne, Hueihan Jhuang, Estibaliz Garrote, Tomaso A. Poggio, and Thomas Serre (2011). “HMDB: A large video database for human motion recognition”. In: *Proceedings of the International Conference on Computer Vision*, pp. 2556–2563.
- M. Pawan Kumar, Philip H. S. Torr, and Andrew Zisserman (2008). “Learning Layered Motion Segmentations of Video”. In: *IJCV*.
- Zihang Lai and Weidi Xie (2019). “Self-supervised Learning for Video Correspondence Flow”. In: *Proceedings of the British Machine Vision Conference*.
- Hsin-Ying Lee, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang (2017). “Unsupervised representation learning by sorting sequences”. In: *Proceedings of the International Conference on Computer Vision*.
- Sungho Lee, Seoung Wug Oh, DaeYeun Won, and Seon Joo Kim (2019). “Copy-and-Paste Networks for Deep Video Inpainting”. In: *Proceedings of the International Conference on Computer Vision*.
- Xiaodan Liang, Lisa Lee, Wei Dai, and Eric P. Xing (2017). “Dual Motion GAN for Future-Flow Embedded Video Prediction”. In: *Proceedings of the International Conference on Computer Vision*.
- Lingjie Liu, Weipeng Xu, Michael Zollhöfer, Hyeonwoo Kim, Florian Bernard, Marc Habermann, Wenping Wang, and Christian Theobalt (2019). “Neural Rendering and Reenactment of Human Actor Videos”. In: *ACM Trans. Graph.* 38.5.
- Stephen Lombardi, Jason Saragih, Tomas Simon, and Yaser Sheikh (2018). “Deep Appearance Models for Face Rendering”. In: *ACM Trans. Graph.* 37.4, 68:1–68:13.
- Xiankai Lu, Wenguan Wang, Jianbing Shen, Yu-Wing Tai, David Crandall, and Steven C. H. Hoi (2020). “Learning Video Object Segmentation from Unlabeled Videos”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.

- Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth (2021). “NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7210–7219.
- Moustafa Meshry, Dan B. Goldman, Sameh Khamis, Hugues Hoppe, Rohit Pandey, Noah Snavely, and Ricardo Martin-Brualla (2019). “Neural Rerendering in the Wild”. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman (2020). “End-to-End Learning of Visual Representations from Uncurated Instructional Videos”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng (2020). “NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis”. In: *Proceedings of the European Conference on Computer Vision*.
- Ishan Misra, C. Lawrence Zitnick, and Martial Hebert (2016). “Shuffle and Learn: Unsupervised Learning using Temporal Order Verification”. In: *Proceedings of the European Conference on Computer Vision*.
- Ajay Nandoriya, Elgharib Mohamed, Changil Kim, Mohamed Hefeeda, and Wojciech Matusik (2017). “Video Reflection Removal Through Spatio-Temporal Optimization”. In: *Proceedings of the International Conference on Computer Vision*.
- Alasdair Newson, Andrés Almansa, Matthieu Fradet, Yann Gousseau, and Patrick Pérez (2014). “Video Inpainting of Complex Scenes”. In: *SIAM Journal on Imaging Sciences, Society for Industrial and Applied Mathematics*.
- Thu H Nguyen-Phuoc, Chuan Li, Stephen Balaban, and Yongliang Yang (2018). “RenderNet: A deep convolutional network for differentiable rendering from 3D shapes”. In: *Advances in Neural Information Processing Systems*.
- Seoung Wug Oh, Sungho Lee, Joon-Young Lee, and Seon Joo Kim (2019). “Onion-Peel Networks for Deep Video Completion”. In: *Proceedings of the International Conference on Computer Vision*.
- Makoto Okabe, Keita Noda, Yoshinori Dobashi, and Ken Anjyo (2020). “Interactive Video Completion”. In: *IEEE Computer Graphics and Applications*.
- Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu (2019). “Semantic Image Synthesis with Spatially-Adaptive Normalization”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Vol. abs/1903.07291.
- Mandela Patrick, Yuki M. Asano, Ruth Fong, João F. Henriques, Geoffrey Zweig, and Andrea Vedaldi (2020). “Multi-modal Self-Supervision from Generalized Data Transformations”. In: *arXiv preprint arXiv:2003.04298*.

- AJ Piergiovanni, Anelia Angelova, and Michael S. Ryoo (2020). “Evolving Losses for Unsupervised Video Representation Learning”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alexander Sorkine-Hornung, and Luc Van Gool (2017). “The 2017 DAVIS Challenge on Video Object Segmentation”. In: *arXiv:1704.00675*.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei (2015). “ImageNet Large Scale Visual Recognition Challenge”. In: *International Journal of Computer Vision* 115.3, pp. 211–252.
- Laura Sevilla-Lara, Deqing Sun, Varun Jampani, and Michael J. Black (2016). “Optical Flow with Semantic Segmentation and Localized Layers”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3889–3898.
- Jonathan Shade, Steven Gortler, Li-wei He, and Richard Szeliski (1998). “Layered depth images”. In: *Proceedings of the 25th annual conference on Computer graphics and interactive techniques*, pp. 231–242.
- Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhöfer (2019). “DeepVoxels: Learning Persistent 3D Feature Embeddings”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein (2019). “Scene Representation Networks: Continuous 3D-Structure-Aware Neural Scene Representations”. In: *Advances in Neural Information Processing Systems*.
- Khurram Soomro, Amir R. Zamir, and Mubarak Shah (2012). “UCF101: A Dataset of 101 Human Actions Classes From Videos in The Wild”. In: *CoRR* abs/1212.0402.
- Pratul P. Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T. Barron (2021). “NeRV: Neural Reflectance and Visibility Fields for Relighting and View Synthesis”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Pratul P. Srinivasan, Richard Tucker, Jonathan T. Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely (2019). “Pushing the Boundaries of View Extrapolation with Multiplane Images”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Deqing Sun, Erik Sudderth, and Michael J. Black (2012). “Layered segmentation and optical flow estimation over time”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, pp. 1768–1775.
- Deqing Sun, Jonas Wulff, Erik Sudderth, Hanspeter Pfister, and Michael J. Black (2013). “A fully-connected layered model of foreground and background flow”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Portland, OR, pp. 2451–2458.

- Tiancheng Sun, Jonathan T. Barron, Yun-Ta Tsai, Zexiang Xu, Xueming Yu, Graham Fyffe, Christoph Rhemann, Jay Busch, Paul Debevec, and Ravi Ramamoorthi (2019). “Single Image Portrait Relighting”. In: *ACM Trans. Graph.* 38.4.
- Zachary Teed and Jia Deng (2020). “RAFT: Recurrent All-Pairs Field Transforms for Optical Flow”. In: *Proceedings of the European Conference on Computer Vision*.
- Ayush Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason M. Saragih, Matthias Nießner, Rohit Pandey, Sean Ryan Fanello, Gordon Wetzstein, Jun-Yan Zhu, Christian Theobalt, Maneesh Agrawala, Eli Shechtman, Dan B. Goldman, and Michael Zollhöfer (2020). “State of the Art on Neural Rendering”. In: *Eurographics*.
- Justus Thies, Michael Zollhöfer, and Matthias Nießner (2019). “Deferred neural rendering: image synthesis using neural textures”. In: *ACM Trans. Graph.* 38.4, 66:1–66:12.
- Justus Thies, Michael Zollhöfer, Christian Theobalt, Marc Stamminger, and Matthias Nießner (2020). “Image-guided Neural Object Rendering”. In: *Proceedings of the International Conference on Learning Representations*.
- Phillip H. S. Torr, Richard Szeliski, and Padmanabhan Anandan (2001). “An integrated Bayesian Approach to Layer Extraction from Image Sequences”. In: *PAMI* 23.3, pp. 297–304.
- Richard Tucker and Noah Snavely (2020). “Single-View View Synthesis with Multiplane Images”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky (2018). “Deep image prior”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9446–9454.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin (2017). “Attention is All You Need”. In: *Advances in Neural Information Processing Systems*.
- Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba (2016a). “Anticipating visual representations from unlabelled video”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- (2016b). “Generating Videos with Scene Dynamics”. In: *Advances in Neural Information Processing Systems*.
- Carl Vondrick, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy (2018). “Tracking Emerges by Colorizing Videos”. In: *Proceedings of the European Conference on Computer Vision*.

- Carl Vondrick and Antonio Torralba (2017). “Generating the Future with Adversarial Transformers”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Chuan Wang, Haibin Huang, Xiaoguang Han, and Jue Wang (2019). “Video Inpainting by Jointly Learning Temporal Structure and Spatial Details”. In:
- Jiangliu Wang, Jianbo Jiao, and Yun-Hui Liu (2020). “Self-supervised Video Representation Learning by Pace Prediction”. In: *Proceedings of the European Conference on Computer Vision*.
- John Y. A. Wang and Edward H. Adelson (1993a). “Layered Representation for Image Sequence Coding”. In: *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*. Vol. 5, pp. 221–224.
- (1993b). “Layered Representation for Motion Analysis”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 361–366.
- (1994). “Representing Moving Images with Layers”. In: *The IEEE Transactions on Image Processing Special Issue: Image Sequence Compression* 3.5, pp. 625–638.
- Tianyu Wang, Xiaowei Hu, Qiong Wang, Pheng-Ann Heng, and Chi-Wing Fu (2020). “Instance Shadow Detection”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro (2018). “Video-to-Video Synthesis”. In: *Advances in Neural Information Processing Systems*.
- Xiaolong Wang and Abhinav Gupta (2015). “Unsupervised Learning of Visual Representations using Videos”. In: *Proceedings of the International Conference on Computer Vision*.
- Xiaolong Wang, Allan Jabri, and Alexei A. Efros (2019). “Learning Correspondence from the Cycle-Consistency of Time”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Donglai Wei, Joseph Lim, Andrew Zisserman, and William T. Freeman (2018). “Learning and Using the Arrow of Time”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Yonatan Wexler, Eli Shechtman, and Michal Irani (2007). “Space-Time Completion of Video”. In: *PAMI*.
- Olivia Wiles, A. Sophia Koepke, and Andrew Zisserman (2018a). “Self-supervised learning of a facial attribute embedding from video”. In: *Proceedings of the British Machine Vision Conference*.
- (2018b). “X2Face: A network for controlling face generation by using images, audio, and pose codes”. In: *Proceedings of the European Conference on Computer Vision*.
- Kar Wai Wong (1994). *Chungking Express*. Jet Tone Production.

- Shangzhe Wu, Tomas Jakab, Christian Rupprecht, and Andrea Vedaldi (2021). “DOVE: Learning Deformable 3D Objects by Watching Videos”. In: *arXiv preprint arXiv:2107.10844*.
- Jonas Wulff and Michael J. Black (2014). “Modeling Blurred Video with Layers”. In: *Proceedings of the European Conference on Computer Vision*. Vol. 8694. Lecture Notes in Computer Science. Springer International Publishing, pp. 236–252.
- (2015). “Efficient Sparse-to-Dense Optical Flow Estimation using a Learned Basis and Layers”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 120–130.
- Dejing Xu, Jun Xiao, Zhou Zhao, Jian Shao, Di Xie, and Yueting Zhuang (2019). “Self-supervised spatiotemporal learning via video clip order prediction”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Jiarui Xu and Xiaolong Wang (2021). “Rethinking Self-Supervised Correspondence Learning: A Video Frame-level Similarity Perspective”. In: *Proceedings of the International Conference on Computer Vision*.
- Ning Xu, Linjie Yang, Yuchen Fan, Jianchao Yang, Dingcheng Yue, Yuchen Liang, Brian Price, Scott Cohen, and Thomas Huang (2018). “YouTube-VOS: Sequence-to-Sequence Video Object Segmentation”. In: *Proceedings of the European Conference on Computer Vision*.
- Rui Xu, Xiaoxiao Li, Bolei Zhou, and Chen Change Loy (2019). “Deep Flow-Guided Video Inpainting”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Zexiang Xu, Kalyan Sunkavalli, Sunil Hadap, and Ravi Ramamoorthi (2018). “Deep Image-based Relighting from Optimal Sparse Samples”. In: *ACM Trans. Graph.* 37.4, 126:1–126:13.
- Tianfan Xue, Michael Rubinstein, Ce Liu, and William T. Freeman (2015). “A Computational Approach for Obstruction-Free Photography”. In: *ACM Transactions on Graphics (Proc. SIGGRAPH)* 34.4.
- Linjie Yang, Yuchen Fan, and Ning Xu (2019). “Video Instance Segmentation”. In: *Proceedings of the International Conference on Computer Vision*.
- Xiuming Zhang, Sean Fanello, Yun-Ta Tsai, Tiancheng Sun, Tianfan Xue, Rohit Pandey, Sergio Orts-Escolano, Philip Davidson, Christoph Rhemann, Paul Debevec, Jonathan T Barron, Ravi Ramamoorthi, and William T Freeman (2021). “Neural Light Transport for Relighting and View Synthesis”. In: *ACM Transactions on Graphics (TOG)* 40.1, pp. 1–17.
- Hang Zhao, Chuang Gan, Wei-Chiu Ma, and Antonio Torralba (2019). “The Sound of Motions”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.

- Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba (2018). “The Sound of Pixels”. In: *The European Conference on Computer Vision (ECCV)*.
- Hao Zhou, Sunil Hadap, Kalyan Sunkavalli, and David Jacobs (2019). “Deep Single Image Portrait Relighting”. In: *Proc. International Conference on Computer Vision (ICCV)*.
- Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely (2018). “Stereo Magnification: Learning View Synthesis using Multiplane Images”. In: *Proceedings of the ACM SIGGRAPH Conference on Computer Graphics*.
- Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A. Efros (2016). “Generative Visual Manipulation on the Natural Image Manifold”. In: *Proceedings of the European Conference on Computer Vision*.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros (2017). “Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks”. In: *Proceedings of the International Conference on Computer Vision*.
- C. Lawrence Zitnick, Sing Bing Kang, Matthew Uyttendaele, Simon Winder, and Richard Szeliski (2004). “High-quality video view interpolation using a layered representation”. In: *ACM Transactions on Graphics*.