

# Motion Segmentation from a Moving Monocular Camera

Yuxiang Huang<sup>1</sup> and John Zelek<sup>1</sup>

**Abstract**—Identifying and segmenting moving objects from a moving monocular camera is difficult when there is unknown camera motion, different types of object motions and complex scene structures. To tackle these challenges, we take advantage of two popular branches of monocular motion segmentation approaches: point trajectory based and optical flow based methods, by synergistically fusing these two highly complementary motion cues at object level. By doing this, we are able to model various complex object motions in different scene structures at once, which has not been achieved by existing methods. We first obtain object-specific point trajectories and optical flow mask for each common object in the video, by leveraging the recent foundational models in object recognition, segmentation and tracking. We then construct two robust affinity matrices representing the pairwise object motion affinities throughout the whole video using epipolar geometry and the motion information provided by optical flow. Finally, co-regularized multi-view spectral clustering is used to fuse the two affinity matrices and obtain the final clustering. Our method shows state-of-the-art performance on the KT3DMoSeg dataset, which contains complex motions and scene structures. Being able to identify moving objects allows us to remove them for map building when using visual SLAM or SFM.

## I. INTRODUCTION

The objective of motion segmentation is to divide a video frame into regions segmented by common motions. Motion segmentation from a moving camera is particularly important as it has various applications in areas like autonomous navigation, robotics and SLAM. In a dynamic scene, the video camera is moving at an unknown velocity with respect to the environment. Such scenarios pose many challenges to motion segmentation methods such as motion degeneracy, motion parallax and motion on the epipolar plane [1]. Existing monocular motion segmentation methods often fail when facing these challenges [2]–[4], or need specific assumptions on the scene structure, object classes or motion types [5]–[7].

In order to overcome these limitations and achieve high quality instance motion segmentation results regardless of scene structures and motion types, we draw our inspiration from two branches of well studied motion segmentation approaches: optical flow based methods and point trajectory based methods. These two types motion cues are not only complementary in nature (long-term vs instantaneous motion), but they can also be used to derive highly complementary geometric and motion models for different motion types and scene structures: points trajectory based methods, when analyzed using epipolar geometry, will fail if the motion is mainly on the epipolar plane, but it is robust to depth variations, perspective effects and motion parallax; on the other hand, optical flow based methods will fail on these challenges, but it is robust to

motions on the epipolar plane. We propose to combine these two complimentary motion cues at object level to obtain a robust and comprehensive motion representation of the scene. By using the state-of-the-art methods for object recognition, detection, segmentation and tracking, we can easily obtain an objectiveness prior (i.e., an initial grouping) of all motion data. This approach enables us to analyze motions of each individual object, which is crucial for both robots and humans to build their situational awareness and scene understanding capabilities [8], [9].

To build our motion segmentation framework, we first leverage state-of-the-art deep learning models [10]–[13] to recognize, detect, segment and track any common objects throughout the video. Then, for every object in the video, we obtain (1) a set of sparse point trajectories for each object and (2) a dense optical flow mask for each object on each frame. By using object-specific sparse point trajectories and optical flow masks, we are able to derive object-specific geometric models (i.e. fundamental matrices based on epipolar geometry) and instantaneous motion models respectively on every frame pair where the object is visible. By fusing these two highly complementary models, we are theoretically able to model the vast majority of motions even in complex scenes. Our experiments show significant improvements on motion segmentation results in challenging scenarios, highlighting the approach’s potential in real-world applications. In summary, the key contributions of our paper are as follows:

- 1) We combine the well-studied fundamental matrix motion model and the optical flow based instantaneous motion model using multi-view spectral clustering to model multiple complex motions in challenging scenes.
- 2) We show how to model different motions at object level by incorporating per-frame objectiveness prior obtained from recent computer vision foundational models.
- 3) We achieve state-of-the-art result on the challenging KT3DMoSeg dataset [4] in terms of both producing high-quality point trajectory clustering and pixel-level masks for individual moving instances.

## II. RELATED WORK

Monocular motion segmentation can be broadly categorized into three groups: (1) Intensity based methods [14]–[18], (2) sparse correspondence based methods [4], [19]–[26] and (3) deep learning methods [5]–[7], [27]–[32].

### A. Intensity Based Methods

Intensity based methods can be further categorized into indirect and direct methods. Indirect methods [15]–[17] rely on pixel-wise correspondences as input, and produce a pixel-wise segmentation mask indicating different motion groups. Such

<sup>1</sup>Yuxiang Huang and <sup>1</sup>John Zelek are with Systems Design Engineering Department, University of Waterloo, Waterloo, N2L 3G1, Canada. {yuxiang.huang, jzelek}@uwaterloo.ca

input is usually optical flow, which assumes the brightness or color intensity of every point in the scene remains the same throughout the whole sequence. In contrast, direct methods [14] do not require optical flow – they combine the two processes of optimizing for the brightness consistency constraint and estimating the motion models together and directly take a pair of images as input. Most recent works on intensity based methods use optical flow based indirect methods, possibly due to the fast advance in optical flow estimation [33], [34]. Such methods usually adopt an iterative optimization approach to estimate the motion model and motion region simultaneously.

Intensity-based methods tend to perform well on scenes without strong depth variations or motion parallax as the motion flow vectors projected to a 2D image from the 3D space are determined by both the depth and the screw motion of objects [35]. Therefore, if the scene contains strong depth variation (e.g. road scenes), intensity-based methods will fail to distinguish if a part of the image is moving independently or is just at a different depth from its surroundings.

#### B. Sparse Correspondence Based Methods

Sparse correspondence based methods can be further categorized into two-frame and multi-frame methods. Two frame methods [19], [20] usually recover motion parameters by solving an iterative energy minimization problem of finding a certain number of geometric models (e.g., fundamental matrices) on a set of matched feature points, to minimize an energy function that evaluates the quality of the overall clustering of correspondences. Multi-frame methods often operate on manually corrected trajectory points obtained from a dense optical flow tracker. Such methods usually perform clustering on affinity matrices constructed using the results of geometric model fitting [4], [25], [26] or pairwise affinities derived from spatio-temporal motion cues and appearance cues [22], [24].

Point trajectories, when analyzed using epipolar geometry, are robust to depth variations in the scene, but are prone to motions on the epipolar plane. [26] uses trifocal tensor as a more robust model to analyze point trajectories. Trifocal tensor is more robust to noise and motions on the epipolar plane, but it is harder to optimize and prone to failure when the three cameras are close to being colinear [1], which can often happen on road scenes. [2], [4] produce promising results by combining multiple geometric models, but still fails to produce a coherent and consistent segmentation on some scenes. Methods using spatio-temporal information and appearance cues [22], [24] suffer from similar issues. In addition, although they tend to perform a bit better than geometric methods on motions with less rigidity [36], they perform worse on scenes with motion parallax or strong camera motions.

#### C. Deep Learning Based Methods

Deep learning based methods usually takes a pair or a sequence of input frames as input and directly produces a either a binary segmentation mask of moving vs static objects [6], [27], [28], or a multi-label segmentation mask showing different objects of different motions [5], [7], [29]–[32]. Most

deep learning methods use CNNs and their network architecture can be broadly summarized to have the following main components: (1) a module to extract the appearance information from consecutive frames, (2) a module to extract motion information from the same frames, (3) a module to fuse the appearance and motion information, and (4) a decoder to generate the final segmentation. These methods are usually fully-supervised and require a large amount of training data. Besides, these methods are often only able to perform binary motion detection (moving vs. static) [7], [28], [31], or they are only able to detect instance motions for specific scenes and limited number of object types they are trained on [5], [30]

### III. METHODOLOGY

#### A. Object Recognition, Detection, Segmentation and Tracking

Figure 1 shows a diagram of our motion segmentation pipeline. In order to identify all motions in a video sequence at object level, we first identify every common object in the video and track their movements throughout the video. We achieve this by using the most recent foundational models in object recognition (Recognize Anything Model) [12], detection (Grounding DINO model) [11] and segmentation (Segment Anything Model) [10], and a state-of-the-art object tracker (DeAOT) [13]. We adapt our preprocessing pipeline from Segment and Track Anything (SAMTrack) [37], which is an object segmentation and tracking framework made of the Grounding DINO model, Segment Anything Model (SAM) and the DeAOT tracker. SAMTrack allows the user to segment and track any specific objects in the video with a text prompt. To make our system fully automatic and universal to all motions in most scenes, we avoid using the user-defined text prompt by adding RAM at the beginning of our pipeline to automatically recognize any common objects in the video. In summary, our whole preprocessing pipeline consists of the following main steps: 1) Use RAM to recognize any common objects in the first frame of the video; 2) Feed the output of RAM as a text prompt to the Grounding DINO model to obtain object bounding boxes; 3) Feed these bounding boxes to SAM to obtain an instance segmentation mask of the first frame. Non-max suppression was used to remove objects with an IoU score  $> 0.5$  and instances whose bounding boxes are larger than half the image are also removed; 4) Use the DeAOT tracker to track each object’s mask throughout the entire video. In order to detect new objects entering the scene in the middle of the video, we run step 1) every  $l$  frames to check if there are new objects. If so, we then run steps 2) to 4) to segment and track the new objects together with the existing objects in the frame. The number  $l$  is video-specific, for example, more dynamic videos with more objects entering and leaving the scene would benefit from a smaller  $l$ .

#### B. Obtaining Point Trajectories and Optical Flow Masks

Once we have the instance segmentation mask for every frame of the video, we then need to obtain a set sparse of point trajectories and a dense optical flow mask on every object at every frame where it’s visible. Essentially, we aim to

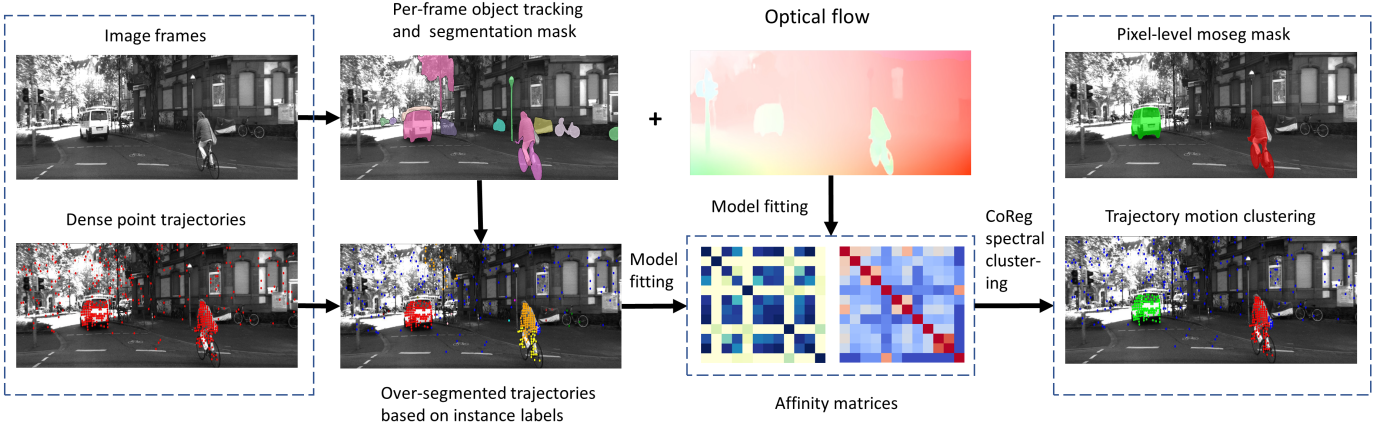


Fig. 1: Motion Segmentation Pipeline

obtain instance point trajectories and instance optical flows. We use one of the state-of-the-art methods [38] to obtain optical flow for all frames, then label each flow vector with the object label of the pixel. Ideally, point trajectories need to be sampled from each object and automatically tracked at each frame using a point tracker (e.g., [39], [40]), however, for benchmarking purposes, we use the manually corrected point trajectories provided by the KT3DMoSeg dataset. We assign every point trajectory with an initial object label using the per-frame instance segmentation masks. If a trajectory does not belong to any object, we label it as the background. Due to inaccuracy in instance segmentation and tracking, a point trajectory can be identified to be on different objects or background in different frames. In such case, we assign its object label to be the most frequent label it is identified as. In the future, we will use a point tracker as well as robust point sampling and occlusion handling techniques to automatically generate point trajectories for all detected objects.

### C. Model Fitting

After obtaining instance point trajectories and optical flow masks for each frame, we compute the motion models of each object to model its motion throughout the video. To compute the epipolar geometry based motion models using point trajectories, we compute a fundamental matrix of each object between every  $f$  frames by solving  $p'TFp = 0$  using the eight-point algorithm. If a degenerate case is encountered for the fundamental matrix, we do not use it. For the optical flow based motion model, we use a full quadratic motion model with 12 parameters to model the instantaneous object screw motion:

$$f(x, y) = (a + bx + cy + dx^2 + exy + fy^2, g + hx + iy + jx^2 + kxy + ly^2) \quad (1)$$

where  $(x, y)$  is the 2D coordinates of the pixels relative to the image center. Since we already have the instance optical flow field, we can obtain the following equation:

$$f(x, y) = (u, v) \quad (2)$$

where  $(u, v)$  is the optical flow vector of the pixel. We fit the function 1 above to the optical flow vectors of every object and solve for the 12 parameters representing the object motion model by optimizing the mean squared error. We use this specific motion model as it's a simplified version of the classic Longuet-Higgins and Pruzdny model equations [41], which model's the instantaneous screw motion of rigid objects at arbitrary depth. Since it's not possible to solve for the depth of each pixel, this motion model assumes the objects' depths are only slightly different. It was shown to perform well on scenes with limited motion parallax [32], nevertheless, it often fails when there is strong motion parallax and depth variations.

### D. Constructing Affinity Matrices

After all fundamental matrices and optical flow motion models are computed, each object will have a fundamental matrix between every  $f$  frames and an optical flow motion model between every two frames. By fitting every object's trajectory points and optical flow vectors to every other object's fundamental matrix and optical flow motion model on the same frame pair, we can obtain the residuals of every object to all other objects' motion models respectively. We use Sampson distance as the residual for fundamental matrix [1] and mean squared error for optical flow motion model. Assuming there are  $k$  objects in the scene, for the  $i$ -th object at the  $m$ -th frame pair, we obtain the following residual vectors under the fundamental matrix and optical flow motion models:

$$\mathbf{r}_{o_i}^m = [r_{o_{i,1}}^m, r_{o_{i,2}}^m, \dots, r_{o_{i,k}}^m],$$

$$\mathbf{r}_{f_i}^m = [r_{f_{i,1}}^m, r_{f_{i,2}}^m, \dots, r_{f_{i,k}}^m]$$

where  $r_{o_{i,k}}^m$  is the mean residual for fitting the parametric motion model of object  $i$  on the optical flow vectors of object  $k$  between frames  $m$  and  $m + 1$ , and  $r_{f_{i,k}}^m$  is the mean Sampson error for fitting the fundamental matrix of object  $i$  on the trajectory points of object  $k$  between frames  $m$  and  $m + f$ . We construct two affinity matrices encapsulating the pairwise motion affinities between each pair of objects using the ordered residual kernel (ORK) [42]. More specifically, for

each object, we sort its residual vectors in ascending order and define a threshold to select the smallest  $t$ -th residual as inliers. We define  $\mathbf{c}_i = \{0, 1\}^K$  as the inlier mask to denote if an object  $i$  is an inlier for each of the  $K$  motion models, and the pairwise motion affinity between objects  $i$  and  $j$  can be computed as  $\mathbf{a}_{ij} = \mathbf{c}_i^T \mathbf{c}_j$ , which denotes the co-occurrence between two objects as an inlier of all motion models. ORK is robust to outliers and makes the affinity matrix more adaptive to different scenes by reducing the need to set scene specific inlier thresholds.

#### E. Co-Regularized Multi-view Spectral Clustering

After constructing the affinity matrices, we use the epsilon-neighborhood scheme [25] to sparsify the affinity matrices. We adapt co-regularized multi-view spectral clustering [43] to fuse the two affinity matrices together to obtain the final clustering of object motions and trajectory points. Co-regularized multi-view spectral clustering uses an regularization term to encourage consensus between different views and is shown to perform well on fusing multiple geometric models for a consistent representation of data [4].

### IV. PRELIMINARY RESULTS & CONCLUSIONS

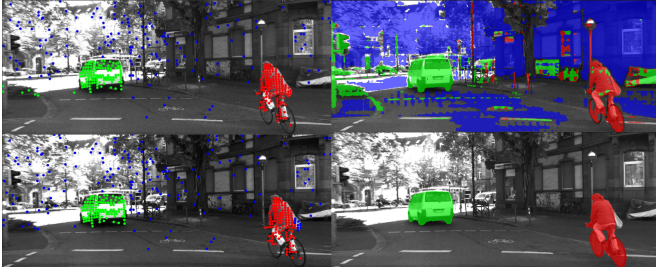


Fig. 2: Trajectory clustering results and the generated segmentation masks by [4] (top) and our method (bottom)

We tested our method on the KT3DMoSeg dataset since it is the only existing dataset involving challenging scenes and multiple complex motions. Since KT3DMoSeg uses pre-defined point trajectories, it could occur that an object segmented by our preprocessing pipeline has fewer than 7 point trajectories. In this case, it's not possible to compute fundamental matrices for the object, but we can still obtain its residual vectors by fitting its trajectory points (if there is any) and the optical flow vectors on the motion models of other objects to compute its pairwise motion affinity scores. Figure 2 shows a qualitative comparison between the segmentation masks generated by our method and a baseline we created using a state-of-the-art method [4] whose code is publicly available. To establish the baseline, we prompt SAM using the clustered trajectories of [4] (top left) to produce the segmentation mask. Even though [4] achieves very low clustering error rate on this sequence, the quality of the generated segmentation mask is still worse than ours since a few wrongly labeled trajectory points can easily mislead SAM. It's also not able to recognize which motion is background. Figure 3 shows qualitative ablation

study and comparison to state-of-the-art methods. Table I shows quantitative results comparing to existing methods. Our method outperforms the state-of-the-art both qualitatively and quantitatively. We also compare the total running time (measured on an Intel 13900K CPU and an NVIDIA RTX 4090 GPU) of our method with [4] for reference. Our method is more than twice as slow since it's not yet optimized.

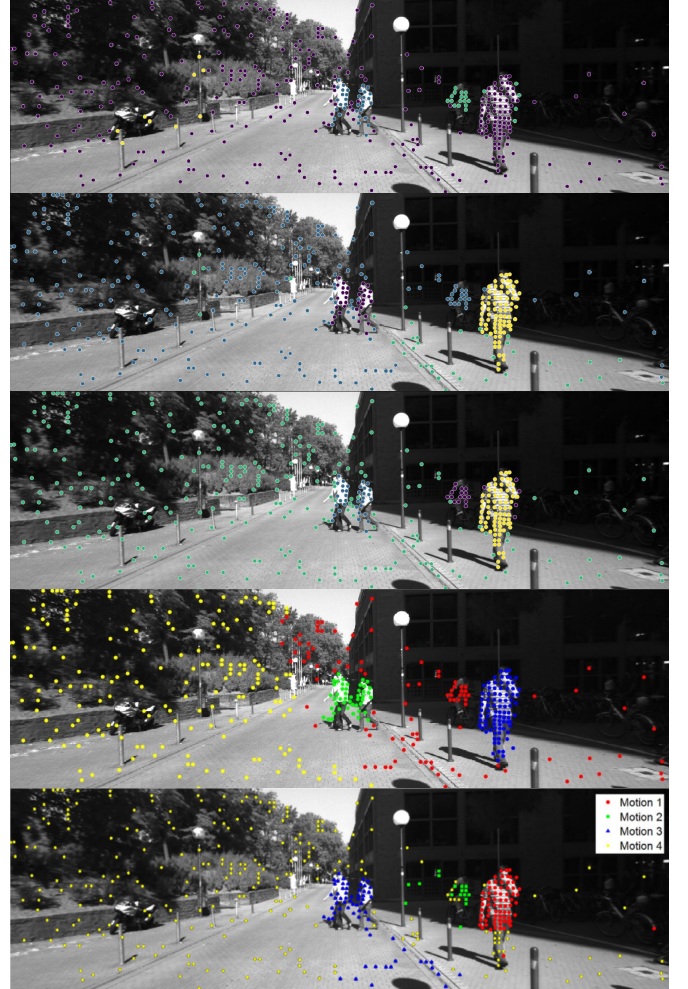


Fig. 3: Qualitative ablation study and comparison with state-of-the-art methods on the Seq38\_Clip02 sequence. From top to bottom: fundamental matrix only, optical flow only, fundamental matrix + optical flow, MVC [4], CMFO [2]

TABLE I: Quantitative results in terms of average classification error (%) and total running time (s). Lower is better

Methods	Avg. Error Rate (%)	Running Time (s)
LSA [44]	38.30	-
GPCA [45]	34.60	-
ALC [46]	24.31	-
SSC [23]	33.88	-
LRR [47]	33.67	-
MVC [4]	10.99	1409.6
CMFO [2]	6.73	-
Ours	<b>5.78</b>	3230.1



## REFERENCES

- [1] R. I. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, 2nd ed. Cambridge University Press, ISBN: 0521540518, 2004.
- [2] Z. Xi, J. Liu, B. Luo, and Q. Qin, "Multi-Motion Segmentation: Combining Geometric Model-Fitting and Optical Flow for RGB Sensors," *IEEE Sensors Journal*, vol. 22, no. 7, pp. 6952–6963, Apr. 2022, conference Name: IEEE Sensors Journal.
- [3] F. Xu, G. Gu, K. Ren, and W. Qian, "Motion Segmentation by New Three-View Constraint from a Moving Camera," *Mathematical Problems in Engineering*, vol. 2015, pp. 1–14, 2015. [Online]. Available: <http://www.hindawi.com/journals/mpe/2015/546580/>
- [4] X. Xu, L. F. Cheong, and Z. Li, "Motion Segmentation by Exploiting Complementary Geometric Models," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. Salt Lake City, UT, USA: IEEE, June 2018, pp. 2859–2867. [Online]. Available: <https://ieeexplore.ieee.org/document/8578400/>
- [5] E. Mohamed, M. Ewaisha, M. Siam, H. Rashed, S. Yogamani, W. Hamdy, M. El-Dakdouky, and A. El-Sallab, "Monocular Instance Motion Segmentation for Autonomous Driving: KITTI InstanceMotSeg Dataset and Multi-Task Baseline," in *2021 IEEE Intelligent Vehicles Symposium (IV)*. Nagoya, Japan: IEEE Press, July 2021, pp. 114–121. [Online]. Available: <https://doi.org/10.1109/IV48863.2021.9575445>
- [6] J. Vertens, A. Valada, and W. Burgard, "SMSnet: Semantic motion segmentation using deep convolutional neural networks," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Sept. 2017, pp. 582–589, iSSN: 2153-0866.
- [7] M. Ramzy, H. Rashed, A. E. Sallab, and S. Yogamani, "RST-MODNet: Real-time Spatio-temporal Moving Object Detection for Autonomous Driving," Dec. 2019, arXiv:1912.00438 [cs, stat] version: 1. [Online]. Available: <http://arxiv.org/abs/1912.00438>
- [8] Y. Jiang, Q. Xu, K. Ma, Z. Yang, X. Cao, and Q. Huang, "What to Select: Pursuing Consistent Motion Segmentation from Multiple Geometric Models," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 2, pp. 1708–1716, May 2021, number: 2. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/16264>
- [9] H. Bavlle, J. L. Sanchez-Lopez, C. Cimorelli, A. Tourani, and H. Voos, "From SLAM to Situational Awareness: Challenges and Survey," *Sensors*, vol. 23, no. 10, p. 4849, Jan. 2023, number: 10 Publisher: Multidisciplinary Digital Publishing Institute. [Online]. Available: <https://www.mdpi.com/1424-8220/23/10/4849>
- [10] F. Rajič, L. Ke, Y.-W. Tai, C.-K. Tang, M. Danelljan, and F. Yu, "Segment Anything Meets Point Tracking," July 2023, arXiv:2307.01197 [cs]. [Online]. Available: <http://arxiv.org/abs/2307.01197>
- [11] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu, and L. Zhang, "Grounding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection," Mar. 2023, arXiv:2303.05499 [cs]. [Online]. Available: <http://arxiv.org/abs/2303.05499>
- [12] Y. Zhang, X. Huang, J. Ma, Z. Li, Z. Luo, Y. Xie, Y. Qin, T. Luo, Y. Li, S. Liu, Y. Guo, and L. Zhang, "Recognize Anything: A Strong Image Tagging Model," June 2023, arXiv:2306.03514 [cs]. [Online]. Available: <http://arxiv.org/abs/2306.03514>
- [13] Z. Yang and Y. Yang, "Decoupling Features in Hierarchical Propagation for Video Object Segmentation," *Advances in Neural Information Processing Systems*, vol. 35, pp. 36324–36336, Dec. 2022. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2022/hash/eb890c36af87e4ca82e8ef7bcb6a284-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/eb890c36af87e4ca82e8ef7bcb6a284-Abstract-Conference.html)
- [14] S. Negahdaripour and B. K. P. Horn, "Direct Passive Navigation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PAMI-9, no. 1, pp. 168–176, Jan. 1987. [Online]. Available: <https://ieeexplore.ieee.org/document/4767884>
- [15] H. Sekkati and A. Mitiche, "A variational method for the recovery of dense 3D structure from motion," *Robotics and Autonomous Systems*, vol. 55, no. 7, pp. 597–607, July 2007. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0921889006001941>
- [16] A. Wedel, D. Cremers, T. Pock, and H. Bischof, "Structure- and motion-adaptive regularization for high accuracy optic flow," in *2009 IEEE 12th International Conference on Computer Vision*, Sept. 2009, pp. 1663–1668, iSSN: 2380-7504.
- [17] P. Bideau and E. Learned-Miller, "It's Moving! A Probabilistic Model for Causal Motion Segmentation in Moving Camera Videos," in *Computer Vision – ECCV 2016*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds. Cham: Springer International Publishing, 2016, vol. 9912, pp. 433–449, series Title: Lecture Notes in Computer Science. [Online]. Available: [http://link.springer.com/10.1007/978-3-319-46484-8\\_26](http://link.springer.com/10.1007/978-3-319-46484-8_26)
- [18] P. Bideau, A. RoyChowdhury, R. R. Menon, and E. Learned-Miller, "The Best of Both Worlds: Combining CNNs and Geometric Constraints for Hierarchical Motion Segmentation," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. Salt Lake City, UT, USA: IEEE, June 2018, pp. 508–517. [Online]. Available: <https://ieeexplore.ieee.org/document/8578158/>
- [19] A. Delong, A. Osokin, H. N. Isack, and Y. Boykov, "Fast approximate energy minimization with label costs," in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, June 2010, pp. 2173–2180, iSSN: 1063-6919.
- [20] H. Isack and Y. Boykov, "Energy-Based Geometric Multi-model Fitting," *International Journal of Computer Vision*, vol. 97, no. 2, pp. 123–147, Apr. 2012. [Online]. Available: <http://link.springer.com/10.1007/s11263-011-0474-7>
- [21] D. Hutchison, T. Kanade, J. Kittler, J. M. Kleinberg, F. Mattern, J. C. Mitchell, M. Naor, O. Nierstrasz, C. Pandu Rangan, B. Steffen, M. Sudan, D. Terzopoulos, D. Tygar, M. Y. Vardi, G. Weikum, T. Brox, and J. Malik, "Object Segmentation by Long Term Analysis of Point Trajectories," in *Computer Vision – ECCV 2010*, K. Daniilidis, P. Maragos, and N. Paragios, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, vol. 6315, pp. 282–295, series Title: Lecture Notes in Computer Science. [Online]. Available: [http://link.springer.com/10.1007/978-3-642-15555-0\\_21](http://link.springer.com/10.1007/978-3-642-15555-0_21)
- [22] T. Brox and J. Malik, "Object segmentation by long term analysis of point trajectories," in *Proceedings of the 11th European conference on Computer vision: Part V*, ser. ECCV'10. Berlin, Heidelberg: Springer-Verlag, Sept. 2010, pp. 282–295.
- [23] E. Elhamifar and R. Vidal, "Sparse Subspace Clustering: Algorithm, Theory, and Applications," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 11, pp. 2765–2781, Nov. 2013. [Online]. Available: <http://ieeexplore.ieee.org/document/6482137/>
- [24] P. Ochs, J. Malik, and T. Brox, "Segmentation of Moving Objects by Long Term Video Analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 6, pp. 1187–1200, June 2014, conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence.
- [25] T. Lai, H. Wang, Y. Yan, T.-J. Chin, and W.-L. Zhao, "Motion Segmentation Via a Sparsity Constraint," *IEEE Transactions on Intelligent Transportation Systems*, vol. 18, no. 4, pp. 973–983, Apr. 2017, conference Name: IEEE Transactions on Intelligent Transportation Systems.
- [26] F. Arrigoni, L. Magri, and T. Pajdla, "On the Usage of the Trifocal Tensor in Motion Segmentation," in *Computer Vision – ECCV 2020*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds. Cham: Springer International Publishing, 2020, vol. 12365, pp. 514–530, series Title: Lecture Notes in Computer Science. [Online]. Available: [https://link.springer.com/10.1007/978-3-030-58565-5\\_31](https://link.springer.com/10.1007/978-3-030-58565-5_31)
- [27] M. Siam, H. Mahgoub, M. Zahran, S. Yogamani, M. Jagersand, and A. El-Sallab, "MODNet: Motion and Appearance based Moving Object Detection Network for Autonomous Driving," in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, Nov. 2018, pp. 2859–2864, iSSN: 2153-0017.
- [28] M. Bosch, "Deep Learning for Robust Motion Segmentation with Non-Static Cameras," Feb. 2021, arXiv:2102.10929 [cs]. [Online]. Available: <http://arxiv.org/abs/2102.10929>
- [29] A. Dave, P. Tokmakov, and D. Ramanan, "Towards Segmenting Anything That Moves," in *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*. Seoul, Korea (South): IEEE, Oct. 2019, pp. 1493–1502. [Online]. Available: <https://ieeexplore.ieee.org/document/9022103/>
- [30] Z. Cao, A. Kar, C. Hane, and J. Malik, "Learning Independent Object Motion From Unlabelled Stereoscopic Videos," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Long Beach, CA, USA: IEEE, June 2019, pp. 5587–5596. [Online]. Available: <https://ieeexplore.ieee.org/document/8953391/>
- [31] M. Faisal, I. Akhter, M. Ali, and R. Hartley, "EpO-Net: Exploiting Geometric Constraints on Dense Trajectories for Motion Saliency," in *2020 IEEE Winter Conference on Applications of Computer Vision (WACV)*. Snowmass Village, CO, USA: IEEE, Mar. 2020, pp. 1873–1882. [Online]. Available: <https://ieeexplore.ieee.org/document/9093589/>

- [32] E. Meunier, A. Badoual, and P. Bouthemy, "EM-Driven Unsupervised Learning for Efficient Motion Segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 4, pp. 4462–4473, Apr. 2023, conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence.
- [33] Z. Teed and J. Deng, "RAFT: Recurrent All-Pairs Field Transforms for Optical Flow," arXiv, Tech. Rep. arXiv:2003.12039, Aug. 2020, arXiv:2003.12039 [cs] type: article. [Online]. Available: <http://arxiv.org/abs/2003.12039>
- [34] D. Sun, C. Herrmann, F. Reda, M. Rubinstein, D. Fleet, and W. T. Freeman, "Disentangling Architecture and Training for Optical Flow," Sept. 2022, arXiv:2203.10712 [cs]. [Online]. Available: <http://arxiv.org/abs/2203.10712>
- [35] A. Mitiche and J. Aggarwal, *Computer Vision Analysis of Image Motion by Variational Methods*, ser. Springer Topics in Signal Processing. Cham: Springer International Publishing, 2014, vol. 10. [Online]. Available: <http://link.springer.com/10.1007/978-3-319-00711-3>
- [36] J. Lezama, K. Alahari, J. Sivic, and I. Laptev, "Track to the future: Spatio-temporal video segmentation with long-range motion cues," in *CVPR 2011*, June 2011, pp. 3369–3376, iSSN: 1063-6919.
- [37] Y. Cheng, L. Li, Y. Xu, X. Li, Z. Yang, W. Wang, and Y. Yang, "Segment and Track Anything," May 2023, arXiv:2305.06558 [cs]. [Online]. Available: <http://arxiv.org/abs/2305.06558>
- [38] D. Sun, C. Herrmann, F. Reda, M. Rubinstein, D. J. Fleet, and W. T. Freeman, "Disentangling Architecture and Training for Optical Flow," in *Computer Vision – ECCV 2022*, ser. Lecture Notes in Computer Science, S. Avidan, G. Brostow, M. Cissé, G. M. Farinella, and T. Hassner, Eds. Cham: Springer Nature Switzerland, 2022, pp. 165–182.
- [39] A. W. Harley, Z. Fang, and K. Fragkiadaki, "Particle Video Revisited: Tracking Through Occlusions Using Point Trajectories," in *Computer Vision – ECCV 2022*, ser. Lecture Notes in Computer Science, S. Avidan, G. Brostow, M. Cissé, G. M. Farinella, and T. Hassner, Eds. Cham: Springer Nature Switzerland, 2022, pp. 59–75.
- [40] C. Doersch, Y. Yang, M. Vecerik, D. Gokay, A. Gupta, Y. Aytar, J. Carreira, and A. Zisserman, "TAPIR: Tracking Any Point with per-frame Initialization and temporal Refinement," June 2023, arXiv:2306.08637 [cs]. [Online]. Available: <http://arxiv.org/abs/2306.08637>
- [41] H. C. Longuet-Higgins and K. Prazdny, "The Interpretation of a Moving Retinal Image," *Proceedings of the Royal Society of London. Series B, Biological Sciences*, vol. 208, no. 1173, pp. 385–397, 1980, publisher: The Royal Society. [Online]. Available: <https://www.jstor.org/stable/35316>
- [42] T.-j. Chin, H. Wang, and D. Suter, "The Ordered Residual Kernel for Robust Motion Subspace Clustering," in *Advances in Neural Information Processing Systems*, vol. 22. Curran Associates, Inc., 2009. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2009/hash/b337e84de8752b27eda3a12363109e80-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2009/hash/b337e84de8752b27eda3a12363109e80-Abstract.html)
- [43] A. Kumar, P. Rai, and H. Daume, "Co-regularized Multi-view Spectral Clustering," in *Advances in Neural Information Processing Systems*, vol. 24. Curran Associates, Inc., 2011. [Online]. Available: [https://papers.nips.cc/paper\\_files/paper/2011/hash/31839b036f63806cba3f47b93af8ccb5-Abstract.html](https://papers.nips.cc/paper_files/paper/2011/hash/31839b036f63806cba3f47b93af8ccb5-Abstract.html)
- [44] J. Yan and M. Pollefeys, "A General Framework for Motion Segmentation: Independent, Articulated, Rigid, Non-rigid, Degenerate and Non-degenerate," in *Computer Vision – ECCV 2006*, ser. Lecture Notes in Computer Science, A. Leonardis, H. Bischof, and A. Pinz, Eds. Berlin, Heidelberg: Springer, 2006, pp. 94–106.
- [45] R. Vidal, R. Tron, and R. Hartley, "Multiframe Motion Segmentation with Missing Data Using PowerFactorization and GPCA," *International Journal of Computer Vision*, vol. 79, no. 1, pp. 85–105, Aug. 2008. [Online]. Available: <https://doi.org/10.1007/s11263-007-0099-z>
- [46] S. Rao, R. Tron, R. Vidal, and Y. Ma, "Motion Segmentation in the Presence of Outlying, Incomplete, or Corrupted Trajectories," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 10, pp. 1832–1845, Oct. 2010, conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence.
- [47] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma, "Robust Recovery of Subspace Structures by Low-Rank Representation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 171–184, Jan. 2013. [Online]. Available: <http://ieeexplore.ieee.org/document/6180173/>