A Deep Moving-camera Background Model

Guy $\text{Erez}^{1[0000-0002-4545-6664]}$, Ron Shapira $\text{Weber}^{1[0000-0003-4579-0678]}$, and Oren Freifeld $^{1[0000-0001-9816-9709]}$

Ben-Gurion University of the Negev, Be'er Sheva, Israel {ergu,ronsha}@post.bgu.ac.il, orenfr@cs.bgu.ac.il

Abstract. In video analysis, background models have many applications such as background/foreground separation, change detection, anomaly detection, tracking, and more. However, while learning such a model in a video captured by a static camera is a fairly-solved task, in the case of a Moving-camera Background Model (MCBM), the success has been far more modest due to algorithmic and scalability challenges that arise due to the camera motion. Thus, existing MCBMs are limited in their scope and their supported camera-motion types. These hurdles also impeded the employment, in this unsupervised task, of end-to-end solutions based on deep learning (DL). Moreover, existing MCBMs usually model the background either on the domain of a typically-large panoramic image or in an online fashion. Unfortunately, the former creates several problems, including poor scalability, while the latter prevents the recognition and leveraging of cases where the camera revisits previously-seen parts of the scene. This paper proposes a new method, called DeepMCBM, that eliminates all the aforementioned issues and achieves state-of-the-art results. Concretely, first we identify the difficulties associated with joint alignment of video frames in general and in a DL setting in particular. Next, we propose a new strategy for joint alignment that lets us use a spatial transformer net with neither a regularization nor any form of specialized (and non-differentiable) initialization. Coupled with an autoencoder conditioned on unwarped robust central moments (obtained from the joint alignment), this yields an end-to-end regularization-free MCBM that supports a broad range of camera motions and scales gracefully. We demonstrate DeepMCBM's utility on a variety of videos, including ones beyond the scope of other methods. Our code is available at https://github.com/BGU-CS-VIL/DeepMCBM.

Keywords: unsupervised; background model; background subtraction; moving camera; joint alignment; regularization-free; deep learning; video analysis.

1 Introduction

The unsupervised video-analysis task this paper focuses on is learning a background model in a video captured by a moving camera. In the simpler case where the camera is static, such models have been used successfully in many computer-vision applications such as background/foreground separation, change or anomaly detection, and tracking. Static-camera solutions, however, cannot be easily extended to the moving-camera case since we do not know, a-priori, how the video frames should be aligned to each



(a) Examples for several input frames



(b) Alignment, visualized via the mean panoramic image (computed from the entire video)



(c) Background estimation using the Conditional Autoencoder

Fig. 1: Typical results of the proposed module. Note that despite the fact that the dog spent long times being static in two locations (as is evident by the corresponding ghosting effects in (b)) the model succeeded in eliminating it from the background.

other. Thus, most of the tools traditionally used in background models become less applicable; *e.g.*, methods based on learning a low-dimensional subspace via Robust Principal Component Analysis (RPCA) assume that the frames are aligned to each other.

Seemingly, there is a straightforward solution: "simply" align the frames to each other to reduce the problem back to the static-camera case, and then build a static-camera background model based on the aligned frames. However, this is more complicated than it might seem. First, the alignment problem itself is often difficult. For example, methods based on creating a panoramic image by sequentially aligning each pair of consecutive frames suffer from drift errors. Moreover, such methods cannot exploit the information conveyed in situations where the camera revisits (possibly from a different viewpoint) a previously-seen region in the scene. This, among other considerations, motivates solutions based on Joint Alignment (JA) of the frames. However, even in this formulation the problem is often still hard to solve, partially due to reasons we analyze later in § 4. Second, and regardless of how the alignment is done, there is the issue of scalability which pertains to not only the alignment problem itself but also the subsequent learning of the background model: when the accumulative motion of the camera throughout the video is substantial, the domain of the panoramic image can be huge so background models learned in that domain must scale gracefully. Furthermore, in such cases, when a frame is warped (i.e., aligned) towards the panorama, it captures only a small portion of the latter. This means that most of the data in the panoramic version of the warped images is missing. This is problematic in our context since existing

solutions for subspace learning in the presence of missing data usually struggle in such cases. Therefore, the missing-data issue, together with the scalability requirement, considerably complicates the task. Due to the above reasons, the success in the case of a **Moving-camera Background Model (MCBM)** is lagging far behind its static-camera counterpart. Moreover, the difficulties above have also largely prevented the use of Deep Learning (DL) for this task. This is unfortunate not only because the idea of harnessing the power of DL is attractive but also since it hinders the usage of MCBMs within larger end-to-end pipelines.

With this in mind, the goal of this paper is to provide an effective and scalable DL-based MCBM. To that aim, we start by identifying more precisely what makes JA of video frames challenging: first in the general case and then in the more specific DL context. Next, we design a new JA strategy based on a regularization-free Spatial Transformer Net (STN) and a JA loss involving a memory aspect. Our method requires no auxiliary tools (such as the brittle and non-differentiable initialization used in [10]) that would prohibit its usage within end-to-end pipelines. We also propose a new deep module for learning a background model. The model, based on a Conditional Autoencoder (CAE) and the output of the JA module, is learned in the small domain of the input frames instead of the much-larger panoramic domain. This eliminates scalability issues and targets the goal of estimating the background more directly. Importantly, this module too can be used within end-to-end pipelines. Figure 1 demonstrates the type of results obtained by the proposed modules. Taken together, the proposed two modules give rise to a new and highly-effective MCBM method, coined **DeepMCBM**, which supports a broad range of camera motions and scales gracefully. We demonstrate DeepMCBM's utility on a variety of videos, including ones beyond the scope of competing methods.

Our key contributions are: 1) a DL module, for jointly aligning video frames, that relies on an STN-based optimization and a new training strategy that requires neither regularization nor initialization; 2) a DL background-modeling module that leverages the JA via a CAE conditioned on unwarped robust central moments derived from the JA; 3) together, these two modules form an end-to-end unsupervised MCBM that achieves SOTA results, that scales gracefully, and that supports a wide range of camera motions.

2 Related Work

STN [26] is a DL module that learns and applies a parameterized input-dependent spatial transformation. Given a parameterized transformation family and an input image f, the STN's output consists of a parameter vector θ and a warped image obtained by warping f using T^{θ} (a transformation parameterized by θ). During training, the differentiation of a loss propagates through the STN. In practice, however, and despite their elegance, potential strength, and usage in numerous papers, STNs are often hard to train. Part of our solution addresses exactly such a case, where we take an STN-based optimization problem that was thought to be too difficult [10] and show how it can, in fact, be solved easily, without resorting to a regularization or a sophisticated limiting initialization.

Static-camera background models. Early methods were pixelwise (e.g., [41]) but later the focus has shifted to subspace estimation using Robust Principal Component Analysis and its variants (e.g., [44,7,50,20]). While those models usually do not scale

well, there also exist scalable RPCA models (e.g., [21,8]).

Image alignment. In [13,34], pairwise homographies are estimated between consecutive frames while [27] uses a multi-layer homography. An adaptive panoramic image is built in [47,32] while [43] relies on the assumption that a PTZ camera is used. Most of the works above make stringent assumptions about the camera motion and estimate transformations between pairs of images, sometimes even sequentially. This approach, however, can lead to accumulative errors and/or significant distortions. To avoid such issues, AutoStitch [6] employs bundle adjustment. However, publicly-available implementations of AutoStitch scale poorly with the number of images (*e.g.*, cannot handle more than a few hundreds of frames). This is unlike the proposed approach which scales gracefully. Alignment methods relying on depth or expensive 3D information/reconstruction include [35,30,29,46]. Unlike those works, and similarly to, *e.g.*, [10], the JA approach in this paper is purely 2D-based.

MCBMs. Online RPCA methods (*e.g.*, [3,22,19]) were extended to the case of camera jitter [23] as well as more significant motions [18]. DECOLOR [49] is another MCBM, based on motion detection, that is restricted to small motions. IncPCP-PTI [9] targets a PTZ-camera setting by updating a low-dimensional subspace with the help of an estimated rigid motion between consecutive frames. Several MCBMs are built by first aligning the frames to each other, and then, in the usually-large domain of the obtained panoramic image, learning a background model from the warped images using a static-camera background model that can handle missing data (since each warped image covers only a portion of the panoramic domain). A prime example for such methods is PRPCA [34]. Also of note are methods targeting **moving-object detection in a moving camera**; *e.g.*, [48,39,4]. These works, however, cannot detect changes unrelated to motion and also do not scale well.

STN-based JA. As we explain in § 4, STN-based JA poses several difficulties. On that note, the closest work to ours is JA-POLS [10] which handles some of the difficulties via the usage of a non-differentiable and non-robust initialization, together with a fairly-restrictive regularization. While JA-POLS is effective in cases where it is applicable, it is limited in the camera-motion types it supports and is not an end-to-end solution. We will return to JA-POLS in more detail later on.

Learning background models in the panoramic domain. Once alignment is obtained, in principle a background model can be learned. However, panoramic-size models (*e.g.*, [34]) do not scale while using an ensemble of Partially-overlapping Local Subspaces (POLS) [10] is cumbersome and also suffers from the fact the number of models grows with the size of the panorama. Either way, the existing methods do not offer an end-to-end solution that can be used easily within DL pipelines.

3 Preliminaries: Joint Alignment (JA)

Let $(f^n)_{n=1}^N$ be the frames of the input video and assume the size of each frame is $h \times w$ pixels. Let C be the number of input channels; e.g., C=3 for RGB images (the case considered in this paper). Let $\Omega \subset \mathbb{R}^2$ denote the rectangular $h \times w$ common domain of each f^n , and let θ_n denote the (latent) parameter vector of the spatial transformation associated with the sought-after alignment of $f^n: \Omega \to \mathbb{R}^C$. The transformation itself,

denoted by T^{θ_n} , is viewed as an $\mathbb{R}^2 \to \mathbb{R}^2$ map (not just $\Omega \to \mathbb{R}^2$). The value of $d = \dim(\theta_n)$ depends on the transformation family; *e.g.*, in the affine case, d = 6. The warped version of Ω is $\Omega^n \triangleq T^{\theta_n}(\Omega) \triangleq \{x : \exists x' \in \Omega \text{ s.t. } T^{\theta_n}(x') = x\} \subset \mathbb{R}^2$. Mathematically, we define the warped image as $g^n : \Omega^n \to \mathbb{R}^C$ using the equality

$$g^{n}(T^{\boldsymbol{\theta}_{n}}(\boldsymbol{x}')) = f^{n}(\boldsymbol{x}') \quad \forall \boldsymbol{x}' \in \Omega.$$
(1)

However, due to technical reasons related to image warping [42], it is more convenient and customary to define g^n via the inverse transformation of T^{θ_n} :

$$g_{\boldsymbol{x}}^n \triangleq g^n(\boldsymbol{x}) = f^n((T^{\boldsymbol{\theta}_n})^{-1}(\boldsymbol{x})) \quad \forall \boldsymbol{x} \in \Omega^n.$$
 (2)

Note that g^n depends on θ_n and f^n . Let H and W be the height and width, respectively, of a rectangle, denoted by $\Omega_{\text{scene}} \subset \mathbb{R}^2$, that is large enough to contain $\bigcup_n \Omega^n$. We now define a mask that will be useful for reasons to become clear shortly. Let M^Ω be a single-channel $h \times w$ image whose domain is Ω and whose values are all equal to 1. Let $M^n: \Omega_{\text{scene}} \mapsto [0,1]$ be a non-binary $H \times W$ mask obtained by image warping of M^Ω , according to T^{θ_n} , using zero padding and a bilinear interpolation kernel. That is, for any integral location x in Ω_{scene} , the value of M^n at x, denoted by M^n_x , is given by

$$M_{\boldsymbol{x}}^n = \widetilde{M^{\Omega}}_{\boldsymbol{x}'} \qquad \boldsymbol{x}' = T^{-\boldsymbol{\theta}_n}(\boldsymbol{x}) \in \mathbb{R}^2$$
 (3)

where $\widetilde{M^{\Omega}}_{x'}$ is interpolated from the values of M^{Ω} at the 4 integral locations nearest to x' where whenever any of those integral locations falls outside Ω the value of M^{Ω} at that location is taken to be zero. Thus, $M_{x}^{n}=0$ if all those 4 locations are outside Ω , $M_{x}^{n}=1$ it they all fall inside it, and $0< M_{x}^{n}<1$ otherwise. Let $g_{x,c}^{n}$ denote the value of g_{x}^{n} at channel c. We will refer to $p_{x,c}\triangleq (g_{x,c}^{n})_{n=1}^{N}$ where $c\in\{1,\ldots,C\}$ as the C pixel stacks at location x. Similarly, we define the mask stack at location x as $m_{x}\triangleq (M_{x}^{n})_{n=1}^{N}$. Note that $p_{x,c}$ and m_{x} depend on $(\theta_{n})_{n=1}^{N}$. A **joint-alignment loss**, to be minimized w.r.t. $(\theta_{n})_{n=1}^{N}$, may be formulated in terms of

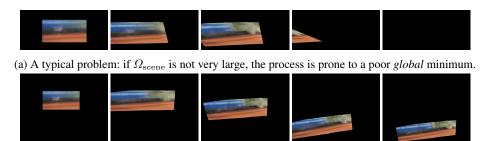
$$\mathcal{L}_{JA} = \operatorname{func}(((p_{\boldsymbol{x},c})_{c=1}^C, m_{\boldsymbol{x}})_{\boldsymbol{x} \in \Omega_{\text{scene}}}). \tag{4}$$

For example, in the early works on *congealing* (e.g., [33,31,25,24]) that loss was based on entropy minimization. Later, other researchers [11,12] showed the benefits of a loss based on least squares. A robust variant (used in [10]) of the latter is

$$\mathcal{L}_{JA} = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{C} \sum_{c=1}^{C} \frac{\sum_{\boldsymbol{x} \in \Omega_{\text{scene}}} M_{\boldsymbol{x}}^{n} \rho_{\text{JA}}(g_{\boldsymbol{x},c}^{n} - \mu_{\boldsymbol{x},c})}{\sum_{\boldsymbol{x} \in \Omega_{\text{scene}}} M_{\boldsymbol{x}}^{n}}$$
(5)

where $\mu_{x,c} = \frac{\sum_{n=1}^{N} M_{x}^{n} g_{x,c}^{n}}{\sum_{n=1}^{N} M_{x}^{n}}$ and $\rho_{\rm JA}$ is a differentiable robust error function [5].

Let μ be the mean of the warped images; *i.e.*, the value of μ at location x and channel c is $\mu_{x,c}$. Note that μ may be viewed as the "moving target" to which the frames should be aligned. It "moves", during the optimization, in the following sense. As the alignment of the frames keeps changing, μ changes too since it is computed using the (weighted) average of the warped images. Assuming that the parameterization $\theta_n \to T^{\theta_n}$ is differentiable and that the transformation family is sufficiently well-behaved (as is the case, e.g.,



(b) A typical problem: drastic spatial changes in μ (note also that the end result is quite blurry).

Fig. 2: Typical problems in JA. Rightmost images are post-convergence results.

with the affine group or, more generally, spaces of diffeomorphisms [15,16,40,38,28]), the loss in Eq. (5) is differentiable. Thus, if θ_n is predicted using an STN (so, in particular, θ_n is a differentiable function of f^n , the STN's input), the loss can, at least in principle, be minimized using standard DL training.

4 Identifying Key Challenges in Solving Joint-alignment Problems

Below we discuss three issues that might arise when solving JA problems: 1) poor *global* minima; 2) the need of regularization; 3) the need of a good initialization.

Usually when trying to minimize a loss, reaching a global minimum is hard or even impossible, and if this feat happens to be achieved, it is deemed to be the ultimate success. Sadly, global minima of $\mathcal{L}_{\mathrm{JA}}$, while being (very) easy to achieve, reflect, in fact, an ultimate failure; e.g., the non-negative $\mathcal{L}_{\mathrm{JA}}$ can attain its global minimum (i.e., zero) when all the frames are shrunk to an infinitesimally-small point. A similar phenomenon occurs if all the frames are warped outside Ω_{scene} (e.g., see Figure 2a) or if the frames are warped such that there will be no pairwise overlap between them.

A popular solution in such cases is adding some type of regularization over $(\theta_n)_{n=1}^N$. However, while various forms of regularization have been suggested, each of them imposes a certain bias; e.g., the regularization term in [31] favors symmetric distributions while the one in [10] pushes the (affine) transformations towards the Special Euclidean group, denoted by SE(2). The implied assumptions in both these cases are limiting. Likewise, penalizing the size of the transformations (e.g.), by penalizing some norm of θ_n) is problematic when the accumulative motion of the camera is large, while regularization favoring temporal smoothness is not always compatible with real camera motions. Another issue is the need of hyperparameter tuning for the weight of the regularization term. Moreover, finding a combination of a regularization type and a weight that will work well for a sufficiently-large variety of videos is difficult.

JA is usually a difficult non-convex problem. Thus, a good initialization can be useful; e.g., in JA-POLS [10] an STN-based JA module had to rely on an initialization based on SE-Sync [37]. The latter provides a useful globally-optimal solution to a different-but-related problem: the estimation of *absolute* transformations that are consistent as possible with noisy measurements of pairwise *relative* transformations between pairs of

frames, where both the latent absolute transformations and the observed relative ones are in SE(2). With that initialization, the STN needs to solve an easier problem and does so over the more expressive Affine group.

There are, however, several problems with the JA approach in [10] (we will later also discuss problems related to the background-modeling approach in [10]). First, preprocessing and heuristics are needed for extracting the relative transformations. Second, in cases where some of the true latent absolute transformations are far from SE(2) (e.g.: when the video contains a significant accumulative variation in the distance between the camera and the scene; when the camera zoom is changing; when there is a strong perspective effect; etc.), the initialization breaks and this leads in turn to JA-POLS' failure. Moreover, SE-Sync is neither robust nor differentiable w.r.t. the input frames. As there is no easy way to differentiate SE-Sync w.r.t. the input frames, the STN-based JA module in JA-POLS cannot be used in an end-to-end DL pipeline.

4.1 An Additional Challenge with Joint Alignment When Using Batches

Typically, due to the data size and as it is almost always the case in DL, the optimization is done batch by batch where each batch consists of a subset (selected at random) of the frames from the entire video. A single epoch then represents a full pass over the entire data, and the frames are reshuffled between epochs. This typically-necessary batch-by-batch processing creates an optimization difficulty which might appear to be minor but is, in fact, far more critical than it may seem (we will revisit this point in § 5.1). The issue is that the mean image μ (from Eq. (5)) is a function of the entire video, not just the frames in the current batch. A seemingly-obvious solution is to hold μ fixed during each epoch – so it does not affect the computation of the loss' gradient – and then, at the end of each epoch, recompute μ . However, a problem that arises with that approach is that the difference between the alignment targets (that is, the previous μ and the recomputed one) in each pair of consecutive epochs might be large, making the optimization difficult since the optimal transformations for one target might be quite far from those that are optimal for the next target. For an illustration, see Figure 2b. A different approach, used in [10], picks the target μ to be the mean of only the (warped) frames in the current batch. Besides the fact that this is somewhat inconsistent with the cost-function formulation, that approach too can cause significant changes in the targets between consecutive batches. The jumping-target problem complicates the optimization more than one may expect. This is especially an issue at the beginning of the process when the frames are completely misaligned. For example, in retrospect, this is partly why JA-POLS [10] had to rely on the SE-Sync-based initialization scheme: as shown in [10], except in the simple case where the accumulative camera motion is small, without that initialization JA-POLS usually fails.

5 The Proposed Method: DeepMCBM

The proposed modules of joint alignment (using an STN) and background modeling (using a CAE) are presented in § 5.1 and § 5.2, respectively. Together, they form the proposed method, DeepMCBM. The goal of the STN straining is 1) to jointly align the

Algorithm 1: Training an STN for Joint Alignment

```
Input: N_{\rm epochs}, N_{\rm batches}, \rho(\cdot), data_loader

Data: (f^n)_{n=1}^N
Output: A trained STN for Joint Alignment

1 Initialize accumulators \mathcal{G} \in \mathbb{R}^{H \times W \times C} and \mathcal{M} \in \mathbb{R}^{H \times W} // see text

2 for e \in \{1, \dots, N_{\rm epochs}\} do

4 (f^b)_{b \in B} \leftarrow \text{data_loader} // Load batch: B \subset \{1, \dots, N\}

5 (\theta_b, g^b)_{b \in B} \leftarrow \text{STN}((f^b)_{b \in B}) // Note that g^b = f^b \circ T^{\theta_b}

6 (M^b)_{b \in B} \leftarrow (M^\Omega \circ T^{\theta_b})_{b \in B} // Warp masks

7 \mathcal{G}, \mathcal{M}, \mathcal{L}_{\text{batch}} \leftarrow \text{Algorithm 2}(\mathcal{G}, \mathcal{M}, (g^b)_{b \in B}, (M^b)_{b \in B}) // Update \mathcal{G} and \mathcal{M}; measure \mathcal{L}_{\text{batch}} (i.e., the batch loss)

8 Perform an optimization step to minimize the \mathcal{L}_{\text{batch}} loss.

9 (\mathcal{G}, \mathcal{M}) \leftarrow (\lambda \mathcal{G}, \lambda \mathcal{M}) // Keep the history, but downweight it
```

video frames, implicitly forming a panoramic image, and 2) to learn how to warp an input frame towards that panoramic image. The goal of the CAE training is to learn the variability in the differences between the panoramic image and the input frames, while taking the warping into account but ignoring the foreground objects. The conditioning is done using the robust version of the panoramic pixelwise mean and variance.

5.1 A Regularization-free Strategy for Joint Alignment

Having identified, in § 4.1, that the jumps in the values of μ cause a major difficulty in the STN-based optimization of \mathcal{L}_{JA} (Eq. (5)), we design a simple but surprisingly-effective optimization strategy, summarized in Algorithm 1 (which, in turn, uses Algorithm 2 as its subroutine). During the training epochs, instead of computing μ using only the current batch (as was done in [10]), or instead of recomputing μ from scratch each epoch, we construct our μ from the warped frames in the current batch while also taking into account, albeit with a lower weight, all the warped frames from the previous epochs as well as the previous batches in the current epoch. The proposed algorithm uses accumulators, denoted by \mathcal{G} and \mathcal{M} . The former is used to accumulate weighted sums of the values of the pixels in the warped frames while the latter serves a similar purpose with the values of the pixels in the warped masks. Concretely, let e denote the index of the current epoch and let e' denote the index of some previous epoch. When evaluating the loss in a batch during epoch e, the contribution of the results from epoch e' becomes smaller and smaller as the "time" difference, e-e', grows. This is done in line 9 in Algorithm 1 by multiplying the accumulators of the warped frames and the warped masks by a positive factor λ where $\lambda < 1$ (we use $\lambda = 0.9$).

As shown in Figure 3, the resulting targets (*i.e.*, the μ sequence formed during the optimization) change smoothly between epochs. Importantly, this behaviour has a profound and fourfold positive effect: **1. No complicated initialization is needed**. As the optimization becomes much easier, the initial transformations are simply taken to be the identity. **2. Regularization-free JA**. No form of regularization on $(\theta_n)_{n=1}^N$ is

Algorithm 2: Update $(\mu, \mathcal{G}, \mathcal{M})$ and measure the loss on the batch

```
Input: \mathcal{G}, \mathcal{M}, (g^b)_{b\in B}, (M^b)_{b\in B}

Output: \mathcal{G}, \mathcal{M}, \mathcal{L}_{\mathrm{batch}}

1 \mathcal{G} \leftarrow \mathcal{G} + \sum_{b=1}^B g^b // update warped-image accumulator

2 \mathcal{M} \leftarrow \mathcal{M} + \sum_{b=1}^B M^b // update warped-mask accumulator

3 \mu \leftarrow \mathbf{0}_{H \times W \times C}

4 for \mathbf{x} \in \{\mathbf{x} : \mathbf{x} \in \Omega_{\mathrm{scene}} \text{ and } \mathcal{M}_{\mathbf{x}} \geq 0\} do in parallel

5 for c \in \{1, \dots, C\} do in parallel

6 \mu_{\mathbf{x},c} \leftarrow \frac{\mathcal{G}_{\mathbf{x},c}}{\mathcal{M}_{\mathbf{x}}}

7 \mathcal{L}_{\mathrm{batch}} \leftarrow \frac{1}{B} \sum_{b=1}^B \frac{1}{C} \sum_{c=1}^C \left[ \left( \sum_{\mathbf{x} \in \Omega_{\mathrm{scene}}} M_{\mathbf{x}}^b \rho(g_{\mathbf{x},c}^b - \mu_{\mathbf{x},c}) \right) / \left( \sum_{\mathbf{x} \in \Omega_{\mathrm{scene}}} M_{\mathbf{x}}^b \right) \right]
```



(a) Compared with Figure 2a, the process is more stable and successful. Also, even when μ nears the border of Ω_{scene} , it never goes outside it.



(b) Compared with Figure 2b, the drastic jumps are eliminated. Also, with the proposed term the results are less affected by the specified size of Ω_{scene} .

Fig. 3: Results analogous to those in Figure 2 except they were obtained with the proposed memory-based approach. Rightmost images are post-convergence results.

needed; e.g., there is no need to worry about the poor global minima from § 4. Since the optimization is gradient-based and since each epoch lingers in the "history" of the process for many epochs before its effective weight decays to zero (due to the repeating multiplications by $\lambda \in (0,1)$, such undesired cases are eliminated altogether. For instance, as the stack of the original frames overlaid over each other (from the first epoch) contributes to the computation of μ , either shrinking the frames to a point or moving them outside Ω_{scene} will incur a loss. Our regularization-free JA is in sharp contrast to many algorithms including classical works (e.g., [31]) and more recent ones (e.g., [10]). **3. Higher expressiveness.** The formulation lets us increase the expressiveness of the transformation family as needed. For example, JA-POLS is so crucially dependent on its SE-Sync initialization and SE-based regularization, that the affine transformations it predicts are nearly in SE(2) themselves. In contrast, our method can not only predict more general transformations in the Affine group but also use broader transformation families. In our experiments we demonstrate this using the group of homographies but one may also try richer STNs such as those based on diffemorphisms [40,1,14]. 4. Our JA module can be used in end-to-end pipelines. This is unlike not only non-DL methods but also JA-POLS [10] whose non-differentiable initialization prevents its JA

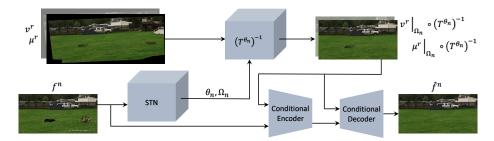


Fig. 4: The background-modeling module. After the STN module was trained using Algorithm 1, the robust panoramic moments, μ^r and v^r , are computed. A CAE is trained for a robust reconstruction task, using the transformation parameters, $(\theta_n)_{n=1}^N$, estimated by the (frozen) STN. The CAE's output is $\widehat{f^n}$, the estimated background associated with f^n and the conditioning is done by (un)warping μ^r and v^r towards each input training image, f^n . During test time the process is similar, except that the transformation being used is the one predicted by the STN.

module from being used in an end-to-end manner. The technical details of the training process appear in our **Supplemental Material (Supmat)**.

5.2 Background Modeling in Ω (not Ω_{scene}) via a Conditional Autoencoder

Upon the training of the STN, the frames become jointly aligned. In principle, at this point all that is left to do is to learn a background model using either non-DL methods (e.g., based on either pixelwise mixture models or RPCA methods; see § 2) or deep ones (such as using a robust loss when training an autoencoder for reconstruction). However, there are several problems with this approach. First, it does not scale well: if the accumulating motion of the camera throughout the video is large, the panoramic image (of the entire scene covered throughout the video) can be huge. Moreover, in such a case even scalable RPCA methods will have to face an additional problem: since the domain of each warped image captures only a small region inside the domain of the panoramic image, it means that most of the pixels will represent missing data. Thus, one would need an RPCA method which can not only scale well but also succeed in situations where more than, say, 90%-95% of the data is missing. Also important is the following. Recall that given an input image, our goal is to estimate a background image, of the same size, that corresponds to that image. Thus, why should we even bother with trying to learn a panoramic-size background model? In [10], the discussion above motivated the learning of multiple local RPCA models and then, for estimating the background of a given image, only a subset of those models whose domains overlapped with the frame of interest were used. That solution, however, means that the number of models to be learned grows with the size of the panoramic image. Moreover, its non-DL formulation was another reason why JA-POLS was not an end-to-end method.

Here we propose a better alternative, whose pipeline is summarized in Figure 4: use a CAE to learn a background model whose domain is small. This has two advantages:

1) It does not compromise the end-to-end nature of the method. 2) We need to learn only a single model (unlike in [10]) and its domain is small, fixed, and does not grow with the size of panoramic image (unlike in PRPCA). Concretely, rather than learning a background model (or models) whose domain is Ω_{scene} , we train a CAE on the original (i.e., non-warped) input frames, using a robust reconstruction error and, for each input frame f^n , conditioning both the encoder and the decoder on (robust versions of) the mean and variance of the pixel stacks, but not before unwarping those central moments from Ω_{scene} towards f^n . We now provide the details. The first and second central moments, denoted by μ^r and v^r , respectively, are C-channel $H \times W$ images defined on Ω_{scence} and computed rubustly using trimmed averaging as follows. Fix $x \in \Omega_{\text{scene}}$, let $N_x = |\{n: M_x^n > 0\}|$, and let $(g_{x,c}^{(1)}, \ldots, g_{x,c}^{(N_x)})$ be the order statistics of $p_{x,c}$. The values of μ^r and v^r at x in channel c are computed, respectively, by

$$\mu_{\boldsymbol{x},c}^{r} = \frac{1}{(1-2\alpha)N_{\boldsymbol{x}}} \sum_{i=\alpha N_{\boldsymbol{x}}}^{(1-\alpha)N_{\boldsymbol{x}}} g_{\boldsymbol{x},c}^{(i)} \qquad v_{\boldsymbol{x},c}^{r} = \frac{1}{(1-2\alpha)N_{\boldsymbol{x}}} \sum_{i=\alpha N_{\boldsymbol{x}}}^{(1-\alpha)N_{\boldsymbol{x}}} (g_{\boldsymbol{x},c}^{(i)} - \mu_{\boldsymbol{x},c}^{r})^{2}.$$
 (6)

Such trimmed averaging is a standard technique for computing robust moments [21]. The trimming parameter, α , was empirically set to $\alpha=0.3$ as it provided a good balance between sample size and robustness. That said, the results when using any other value in the wide range between 20% and almost 50% were similar. Next, when f^n is fed into the CAE, the encoder and the decoder are conditioned by

$$\mu^n \triangleq (\mu^r|_{\Omega_n}) \circ (T^{\theta_n})^{-1} \text{ and } v^n \triangleq (v^r|_{\Omega_n}) \circ (T^{\theta_n})^{-1}$$
 (7)

which are $h \times w$ images (with C channels) defined on Ω and are nothing more than the portion of μ^r and v^r that is relevant for f^n . Using a code whose length was only 4, the CAE was trained with the following loss:

$$\mathcal{L}_{AE} = \sum_{n=1}^{N} \sum_{c=1}^{C} \sum_{\boldsymbol{x}' \in \Omega} \rho_{\text{recon}}(f_{\boldsymbol{x}',c}^{n} - \widehat{f}_{\boldsymbol{x}',c}^{n})$$
(8)

$$\widehat{f}^n = \text{Decoder}(\text{Encoder}(f^n; \mu^n, v^n); \mu^n, v^n)$$
(9)

where \widehat{f}^n is the output of the CAE and $\rho_{\rm recon}$ is a differentiable robust error function. We remark that, by design, the fact that μ^n and v^n are of the same dimensions as the input, f^n , also means it is easy to implement the conditioning via a convolutional layer. For more details about the CAE (whose architecture is based on the AE from [2]) as well as other training details, see our **Supmat**. Finally, $\rho_{\rm recon}$ should usually be more robust than $\rho_{\rm JA}$. The reason is that while in JA the influence of foreground objects is relatively small, in the CAE-based reconstruction it is important, in every pixel, to eliminate the outliers (*i.e.*, the foreground pixels) as much as possible. Thus, we use the smoothed ℓ_1 loss (which is closely-related to Huber's function [5]) for $\rho_{\rm JA}$ and the Geman-McClure error function [17] for $\rho_{\rm recon}$. See **Supmat** for details.

6 Results

We experimented with 4 variants of the proposed DeepMCBM: 1. **Basic/Aff**: This version uses only the STN-based JA module, without the CAE. It estimates the background

by simply unwarping the robust mean towards the input image. The transformations used in the STN belong to the Affine group (the invertibility of the transformations was guaranteed via the matrix exponential; see Supmat). 2. CAE/Aff: This version too uses the Affine STN but also uses the CAE module (for estimating the background). 3. Basic/Hom and 4. CAE/Hom: Similar to Basic/Aff and CAE/Aff, respectively, except that homographies are used instead of affine transformations. We compared those 4 variants with several methods: PRPCA [34]; JA-POLS [10]; PanGAEA [18]; DECOLOR [49]; PCP_PTI [9]; PRAC [19]. The 13 videos that we tested on are ones typically used for evaluation of methods in this area and are taken from well-known datasets [45,36]. Those movies cover camera motions in a variety of types, sizes, speed, zoom changes, etc. It should be noted that, due to their scalability limitaitons, PRPCA and PanGAEA could not run on the ContinuousPan video as the covered scene in the latter was too large. JA-POLS failed running on **zoomInZoomOut** (the significant zoom changes broke its key assumption). ?? contains a visual comparison, on select example videos, of DeepM-CBM (in its CAE/Hom variant), PRPCA, JA-POLS, and PanGAEA. Results of the other (and less successful) methods (DECOLOR; PCP_PTI; PRAC), as well as more visual results (including videos) are in the **Supmat**.

Given an estimate of the background, subtracting it from the original frame yields a difference that can serve to determine foreground/background separation. To quantify the results in a threshold-independent way, for each method and each video we computed the Receiver Operating Characteristic (ROC) curve (using the ground truth) and its Area Under the Curve (AUC). The ROC curves are included in **Supmat**. We emphasize that our method is unsupervised and the ground truth information was used only for evaluation. Table 1, summarizing the AUC results, shows that DeepMCBM, especially with its CAE variants, is, overall, the leading method. In cases where DeepMCBM is not the first it is typically the runner-up. Moreover, unlike some competitors, DeepMCBM was applicable in all cases considered. The visual examples also illustrate how the CAE helps achieving a better estimate of the background than that one obtained by merely using the unwarped μ^r . We remark that our fixed code size, 4, is so small since: 1) the goal is not a typical reconstruction but to filter out foreground objects; 2) our AE is conditional so it is unsurprising a small size suffices. We could have made the code size video-dependent and thus improve results even further, but felt that a fixed size is simpler and makes a comparison with other methods fairer.

Predicting background for previously-unseen misaligned frames. In the comparison above, we focused on background/foreground estimation in the input videos on which the competing models (ours included) were learned. However, like JA-POLS, but unlike all the other methods, our method can predict the background in frames that were not included in the learning (more accurately, some of the competing methods can predict the background in the next constitutive frame, but they are unable to do so for misaligned frames in general such as those that are not consecutive). Due to space limits, we demonstrate that capability of DeepMCBM in the **Supmat**.

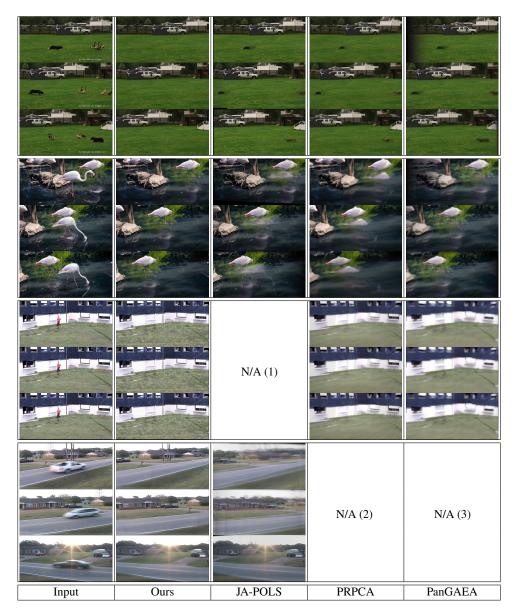


Fig. 5: Visual Comparison: Select Results. Note the ghosting artifacts. N/A (1) large zoom changes failed JA-POLS completely, N/A (2) out-of-memory on a 256GB RAM machine, N/A (3) failed to run: Matlab process was killed

Ablation Study. As Table 1 shows, the AE usually improves performance. In particular, its role is especially important when a foreground object spends a long portion of time in a static position (*e.g.*, the dog in the dog-gooses or the flamingo). In such

	DeepMCBM (Ours)									
Sequence	Basic/Aff	CAE/Aff	Basic/Hom	CAE/Hom	[34]	[18]	[9]	[10]	[19]	[49]
bmx-trees	.898	.896	.916	.908	.894	.786	.837	.930	.664	.737
boxing-fisheye	.924	.893	.927	.898	.935	.932	.728	.892	.627	.763
breakdance-flare	.931	.933	.953	.963	.960	.972	.740	.897	.806	.667
continuousPan	.897	.940	.895	.938	N/A	N/A	.846	.449	.656	.760
dog-gooses	.954	.984	.955	.984	.942	.917	.721	.947	.747	.886
flamingo	.962	.980	.961	.980	.891	.957	.638	.947	.560	.656
horsejump-high	.932	.942	.932	.943	.958	.908	.783	.914	.713	.892
sidewalk	.886	.908	.889	.932	.812	.702	.635	.851	.780	.935
stroller	.877	.885	.740	.756	.762	.904	.594	.807	.613	.721
stunt	.963	.979	.961	.978	.959	.954	.899	.930	.711	.781
swing	.880	.877	.887	.897	.942	.879	.805	.874	.722	.812
tennis	.960	.961	.959	.963	.943	.929	.831	.932	.787	.852
zoomInZoomOut	.981	.994	.981	.994	.979	.958	.720	N/A	.885	.957

Table 1: AUC scores for each method on each sequence.

cases, the robust mean alone still tends to capture some "ghosting" artifacts (as usually do all the competing methods) while the CAE helps correctly identifying that object as belonging to the foreground. The importance of the memory-based approach was also demonstrated in Figure 2 and Figure 3. In particular, the JA failures in Figure 2 imply that no subsequent background model could be built there, making a quantitative comparison (between using the memory term and not using it) a moot point. Finally, note that a basic (*i.e.*, unconditional AE) that knows nothing about the alignment has no chance here as it can only either simply reconstruct the entire frames (*i.e.*, with the undesired foreground objects) or fail in the reconstruction. Thus, when simply dropping the conditioning from our CAE, the resulting AE fails badly in background modeling; *e.g.*, its AUC for the Tennis video is 0.701 while DeepMCBM's AUC score is 0.963.

7 Conclusion

The proposed DeepMCBM is an end-to-end DL solution for modeling background in a video from a moving camera. It supports a wide range of camera-motion types and sizes, scales gracefully, and achieves SOTA results. While we experimented with either affine transformations or homographies, DeepMCBM also supports more expressive transformations. The proposed regularization-free STN-based JA strategy may find usage in other applications, thereby the potential impact of this work may be broader than MCBMs. One limitation of our work is that, since DL involved, the training is slower in comparison to some competitors (JA-POLS excluded). However, we believe the SOTA results together with the other benefits DeepMCBM brings (end-to-end; scalability; the ability to predict background for previously-unseen misaligned frames; *etc.*) justifies it. The main failure case of the method is when foreground objects are large and much closer to the camera than the background is.

Acknowledgements. This work was supported in part by the Lynn and William Frankel Center at BGU CS and by Israel Science Foundation Personal Grant #360/21. G.E. was also funded by the VATAT National excellence scholarship for female Master's students in Hi-Tech-related fields.

References

- Balakrishnan, G., Zhao, A., Sabuncu, M.R., Guttag, J., Dalca, A.V.: An unsupervised learning model for deformable medical image registration. In: CVPR (2018) 9
- Ballé, J., Laparra, V., Simoncelli, E.P.: End-to-end optimized image compression. In: ICLR (2017) 11
- Balzano, L., Nowak, R., Recht, B.: Online identification and tracking of subspaces from highly incomplete information. In: Allerton (2010) 4
- Berger, M., Seversky, L.M.: Subspace tracking under dynamic dimensionality for online background subtraction. In: CVPR (2014) 4
- 5. Black, M.J., Rangarajan, A.: On the unification of line processes, outlier rejection, and robust statistics with applications in early vision. IJCV (1996) 5, 11
- Brown, M., Lowe, D.G.: Automatic panoramic image stitching using invariant features. IJCV (2007) 4
- 7. Candès, E.J., Li, X., Ma, Y., Wright, J.: Robust principal component analysis? JACM (2011) 3
- 8. Chakraborty, R., Hauberg, S., Vemuri, B.C.: Intrinsic grassmann averages for online linear and robust subspace learning. In: CVPR (2017) 4
- 9. Chau, G., Rodríguez, P.: Panning and jitter invariant incremental principal component pursuit for video background modeling. In: ICCV (2017) 4, 12, 14
- 10. Chelly, I., Winter, V., Litvak, D., Rosen, D., Freifeld, O.: JA-POLS: a moving-camera background model via joint alignment and partially-overlapping local subspaces. In: CVPR (2020) 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 14
- 11. Cox, M., Sridharan, S., Lucey, S., Cohn, J.: Least squares congealing for unsupervised alignment of images. In: CVPR (2008) 5
- Cox, M., Sridharan, S., Lucey, S., Cohn, J.: Least-squares congealing for large numbers of images. In: ICCV (2009) 5
- 13. Cuevas, C., Mohedano, R., García, N.: Statistical moving object detection for mobile devices with camera. In: ICCE (2015) 4
- 14. Dalca, A., Rakic, M., Guttag, J., Sabuncu, M.: Learning conditional deformable templates with convolutional networks. In: NeurIPS (2019) 9
- Freifeld, O., Hauberg, S., Batmanghelich, K., Fisher III, J.W.: Highly-expressive spaces of well-behaved transformations: Keeping it simple. In: ICCV (2015) 6
- 16. Freifeld, O., Hauberg, S., Batmanghelich, K., Fisher III, J.W.: Transformations based on continuous piecewise-affine velocity fields. IEEE TPAMI (2017) 6
- 17. Geman, S., McClure, D.E.: Statistical methods for tomographic image reconstruction. In: BISI (1987) 11
- 18. Gilman, K., Balzano, L.: Panoramic video separation with online Grassmannian robust subspace estimation. In: ICCV Workshops (2019) 4, 12, 14
- 19. Guo, H., Qiu, C., Vaswani, N.: Practical reprocs for separating sparse and low-dimensional signal sequences from their sum—part 1. In: ICASSP (2014) 4, 12, 14
- Guyon, C., Bouwmans, T., Zahzah, E.H.: Foreground detection via robust low rank matrix decomposition including spatio-temporal constraint. In: ACCV (2012) 3
- Hauberg, S., Feragen, A., Black, M.J.: Grassmann averages for scalable robust pca. In: CVPR (2014) 4, 11
- 22. He, J., Balzano, L., Szlam, A.: Incremental gradient on the grassmannian for online foreground and background separation in subsampled video. In: CVPR (2012) 4
- 23. He, J., Zhang, D., Balzano, L., Tao, T.: Iterative grassmannian optimization for robust image alignment. Image and Vision Computing (2014) 4

- Huang, G., Mattar, M., Lee, H., Learned-Miller, E.G.: Learning to align from scratch. In: NIPS (2012) 5
- Huang, G.B., Jain, V., Learned-Miller, E.: Unsupervised joint alignment of complex images. In: ICCV (2007) 5
- Jaderberg, M., Simonyan, K., Zisserman, A., et al.: Spatial transformer networks. In: NeurIPS (2015) 3
- 27. Jin, Y., Tao, L., Di, H., Rao, N.I., Xu, G.: Background modeling from a free-moving camera by multi-layer homography algorithm. In: ICIP (2008) 4
- 28. Kaufman, I., Weber, R.S., Freifeld, O.: Cyclic diffeomorphic transformer nets for contour alignment. In: ICIP (2021) 6
- Kendall, A., Grimes, M., Cipolla, R.: Posenet: A convolutional network for real-time 6-dof camera relocalization. In: ICCV (2015) 4
- Klein, G., Murray, D.: Parallel tracking and mapping for small ar workspaces. In: International Symposium on Mixed and Augmented Reality (2007) 4
- 31. Learned-Miller, E.G.: Data driven image models through continuous joint alignment. IEEE TPAMI (2006) 5, 6, 9
- 32. Meneghetti, G., Danelljan, M., Felsberg, M., Nordberg, K.: Image alignment for panorama stitching in sparsely structured environments. In: Scandinavian Conference on Image Analysis (2015) 4
- 33. Miller, E.G., Matsakis, N.E., Viola, P.A.: Learning from one example through shared densities on transforms. In: CVPR (2000) 5
- Moore, B.E., Gao, C., Nadakuditi, R.R.: Panoramic robust PCA for foreground–background separation on noisy, free-motion camera video. IEEE Transactions on Computational Imaging (2019) 4, 12, 14
- 35. Newcombe, R.A., Lovegrove, S.J., Davison, A.J.: Dtam: Dense tracking and mapping in real-time. In: ICCV (2011) 4
- Pont-Tuset, J., Perazzi, F., Caelles, S., Arbeláez, P., Sorkine-Hornung, A., Van Gool, L.: The 2017 davis challenge on video object segmentation. arXiv preprint arXiv:1704.00675 (2017)
- Rosen, D.M., Carlone, L., Bandeira, A.S., Leonard, J.J.: SE-Sync: A certifiably correct algorithm for synchronization over the special Euclidean group. The International Journal of Robotics Research (2019) 6
- 38. Shapira Weber, R., Eyal, M., Skafte Detlefsen, N., Shriki, O., Freifeld, O.: Diffeomorphic temporal alignment nets. In: NeurIPS (2019) 6
- 39. Sheikh, Y., Javed, O., Kanade, T.: Background subtraction for freely moving cameras. In: ICCV (2009) 4
- 40. Skafte Detlefsen, N., Freifeld, O., Hauberg, S.: Deep diffeomorphic transformer networks. In: CVPR (2018) 6, 9
- 41. Stauffer, C., Grimson, W.E.L.: Adaptive background mixture models for real-time tracking. In: CVPR (1999) 3
- Szeliski, R.: Computer vision: algorithms and applications. Springer Science & Business Media (2010) 5
- Thurnhofer-Hemsi, K., López-Rubio, E., Domínguez, E., Luque-Baena, R.M., Molina-Cabello, M.A.: Panoramic background modeling for ptz cameras with competitive learning neural networks. In: IJCNN (2017) 4
- 44. De la Torre, F., Black, M.J.: Robust principal component analysis for computer vision. In: ICCV (2001) 3
- 45. Wang, Y., Jodoin, P.M., Porikli, F., Konrad, J., Benezeth, Y., Ishwar, P.: Cdnet 2014: an expanded change detection benchmark dataset. In: CVPR Workshop (2014) 12
- 46. Wu, C.: Towards linear-time incremental structure from motion. In: International Conference on 3D Vision (2013) 4

- 47. Xue, K., Liu, Y., Chen, J., Li, Q.: Panoramic background model for PTZ camera. In: International Congress on Image and Signal Processing (2010) 4
- 48. Yalcin, H., Hebert, M., Collins, R., Black, M.J.: A flow-based approach to vehicle detection and background mosaicking in airborne video. In: CVPR (2005) 4
- 49. Zhou, X., Yang, C., Yu, W.: Moving object detection by detecting contiguous outliers in the low-rank representation. TPAMI (2012) 4, 12, 14
- 50. Zhou, Z., Li, X., Wright, J., Candes, E., Ma, Y.: Stable principal component pursuit. In: ISIT (2010) 3

An End-to-end Moving Camera Background Model

Supplemental Material

Guy $\text{Erez}^{1[0000-0002-4545-6664]}$, Ron Shapira $\text{Weber}^{1[0000-0003-4579-0678]}$, and Oren Freifeld $^{1[0000-0001-9816-9709]}$

Ben-Gurion University of the Negev, Be'er Sheva, Israel {ergu,ronsha}@post.bgu.ac.il, orenfr@cs.bgu.ac.il

Abstract. The supplemental material includes, in addition to this document, select examples of **videos showing the results of our method next to the original videos**. Additionally, and due to space limits, additional visual comparisons between the different methods (using PDF presentations so it be will easy to browse back and forth between the frames) are available at https://github.com/BGU-CS-VIL/DeepMCBM.

As for this document, it contains the following:

- 1. A comparison of ROC curves of different methods on various videos. These curves correspond to the AUC values reported in Table 1 in the paper.
- 2. A demonstration of DeepMCBM's capability of predicting the background in previously-unseen misaligned frames from the video it was trained on (a capability lacking in most other methods, except JA-POLS [4]).
- 3. The details of the robust error functions that we used.
- 4. The technical details of the evaluation procedures.
- 5. The technical details of the architecture and training process.
- 6. An explanation how, in the affine Spatial Transformer Net, the invertibility of the affine transformations is guaranteed via the matrix exponential.

Erez et al.

1 ROC Curves (Related to the Values from Table 1 in the Paper)

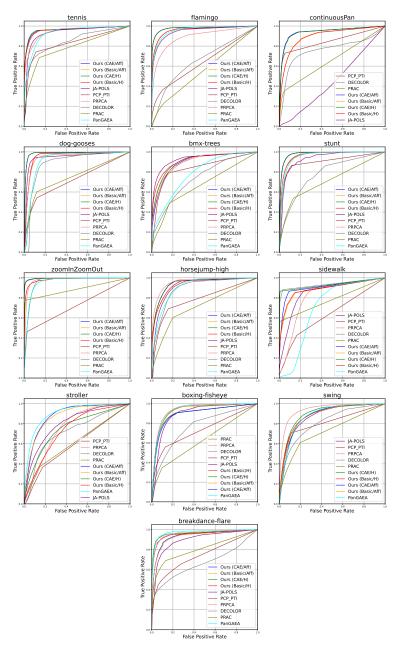


Fig. 1: A comparison of ROC curves of different methods on various videos. Note that sometimes some curves coincide (*e.g.*, our CAE/Aff and CAE/Hom on the flamingo video).

2 Predicting Background to Previously-unseen Misaligned Frames

In the paper, our comparison of DeepMCBM to the other methods focused on their own playground: *i.e.*, given an input sequence of video frames, the task was viewed as an optimization problem whose overarching goal is to estimate the background in those frames. However, like JA-POLS [4] but unlike the other competitors, DeepMCBM, being based on (unsupervised) learning, also possesses a generalization capability in the sense it can estimate the background in previously-unseen misaligned frames from the video it was trained on. In contrast, other methods (JA-POLS excluded), do not have a readily-available mechanism that enables them to receive such misaligned frames and predict their alignment and background estimation – at least not without solving additional optimization steps. Note that in a static camera background model this is a non-issue since, by definition, the frames (both train and test) are always aligned; however, in the moving-camera case the situation is different due to the misalignment. Remark: Note that this generalization capability should not be confused with an online-learning setting – which some competitors aim for – where each time the next consecutive frame arrives their model is being updated (using further optimization).

To showcase the generalization capability, in each of the "tennis" and "flamingo" sequences (from [10]) we partitioned the video sequence into train and test sets. We did it by letting the test set (in each of the two videos) consist of 10% of the frames (chosen at random from the original sequence) while letting the train test consist of the remaining 90% of the frames. In each video we let both DeepMCBM and its competitor, JA-POLS, train only on the train set. Next, we evaluated the models, using their prediction functionalities, on the test sets. Figure 2 and Table 1 summarize the results in terms of the corresponding ROC curve and AUC scores, respectively, and show that DeepMCBM outperforms JA-POLS in this type of evaluation as well. Note that here we intentionally used our CAE/Aff variant (and not CAE/Hom) to highlight the fact that even when using the same transformation type as JA-POLS (*i.e.*, affine), DeepMCBM beats the latter.

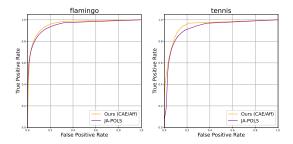


Fig. 2: A typical comparison of ROC curves, between the proposed method and JA-POLS [4], of the performance on test sets.

Sequence	DeepMCBM JA-POLS [4]					
tennis test	0.944	0.923				
tennis train	0.965	0.936				
flamingo test	0.957	0.949				
flamingo train	0.980	0.949				

Table 1: AUC scores for test and train sets. DeepMCBM in its CAE+Affine version

3 The Robust Error Functions

We have used the smoothed ℓ_1 loss (which is closely-related to Huber's function [3]) for $\rho_{\rm JA}$ and Geman-McClure's error function [5] for $\rho_{\rm recon}$:

$$\rho_{\mathrm{JA}}(\varepsilon;\beta) = \begin{cases} 0.5\varepsilon^2/\beta & \text{if} \quad |\varepsilon| \le \beta \\ |\varepsilon| - 0.5\beta & \text{otherwise} \end{cases} \qquad \rho_{\mathrm{recon}}(\varepsilon;s) = \frac{\varepsilon^2}{\varepsilon^2 + s^2} \qquad (1)$$

where in all of our experiments we used $\beta = 0.35$ and s = 0.05.

4 Evaluation Details

In this section we explain how the ROC curves and AUC scores were obtained.

Given an input sequence of N frames, $(f^n)_{n=1}^N$, **ground-truth annotations** of foreground pixels, $(a^n)_{n=1}^N$, and a background estimation sequence, $(\widehat{f}^n)_{n=1}^N$, we compute the pixelwise squared error $(E^n)_{n=1}^N$, where $E_x^n = \frac{1}{C}\sum_{c=1}^C (f_{x,c}^n - \widehat{f}_{x,c}^n)^2$ is the error in pixel $x \in \Omega$ (recall that Ω , a rectangle of height h and width w, was defined in the paper as the common domain of the original input frames), averaged across C channels. To have a unified comparison measure applicable in all videos and for all methods, we scale to the [0,1] interval where the scaling is done (for the method under consideration) w.r.t. the entire video sequence. That is, we define the scaled error of that method in the specific video as

$$\widetilde{E}_{\boldsymbol{x}}^{n} = \frac{E_{\boldsymbol{x}}^{n} - \min_{n, \boldsymbol{x}'}(E_{\boldsymbol{x}'}^{n})}{\max_{n, \boldsymbol{x}'}(E_{\boldsymbol{x}'}^{n}) - \min_{n, \boldsymbol{x}'}(E_{\boldsymbol{x}'}^{n})}.$$
(2)

Next, for each threshold value α (from a discrete set of evenly-spaced points between 0 and 1) we computed the True Positive Rate (TPR) and False Positive Rate (FPR):

$$TPR = \frac{1}{N \cdot h \cdot w} \sum_{n=1}^{N} \sum_{\boldsymbol{x} \in \Omega} \mathbb{1}(a_{\boldsymbol{x}} = 1 \wedge \widetilde{E}_{\boldsymbol{x}}^{n} \ge \alpha);$$
 (3)

$$FPR = \frac{1}{N \cdot h \cdot w} \sum_{n=1}^{N} \sum_{\boldsymbol{x} \in \Omega} \mathbb{1}(a_{\boldsymbol{x}} = 0 \wedge \widetilde{E}_{\boldsymbol{x}}^{n} \ge \alpha).$$
 (4)

6

Computing the TPR and FPR for multiple threshold values lets us obtain the Receiver Operating Characteristic (ROC) curve for this (method, sequence) pair. Figure 1 shows the resulted ROC curves while Area Under the (ROC) Curve (AUC) scores are summarized in Table 1 in the paper.

Architecture and Training Details

DeepMCBM's pipeline consists of two main parts:

- 1. a Spatial Transformer Net (STN);
- 2. a Conditional Autoencoder (CAE).

First, an STN is trained to learn an input-dependent function, that predicts the transformation parameter vector $\boldsymbol{\theta}$. These predictions are used to create the panoramic central moments as well as, later, to unwarp those moments back towards the input image. Once the STN training process has converged, the STN module is frozen and the CAE is trained to learn a background model in the input domain, where the conditioning is on the aforementioned unwapred moments. Below we describe the architecture and training details of the STN and CAE modules.

5.1 The STN

The STN we used consisted of a backbone and tailored regression heads for each available transformation type (in our case, either affine of homographies). The backbone we used was ResNet18 [6] whose output size was 1000. Each regression head consisted of two dense layers of size 1000×32 and $32 \times d$, where d is the dimension of the transformation parameters vector $\boldsymbol{\theta}$ (d=6 for affine transformations and d=8 for homographies). In both the regressor heads we used a ReLU [1] activation function. In all our experiments, we trained the backbone and the Affine head for the first 3000 epochs. In the variants that used homographies, we first did repeat the process above (that is, 3000 epochs for training the backbone and the affine head) and then switched to train the homographic head for additional 3000 epochs (recall that our affine transformations are invertible and that such transformations are a particular case of homographies; thus, the results from the affine head provide a good initialization for the homographic head). Thus, in total, the STN module was trained for about 6000 epochs with a step learning rate scheduler, decreasing the learning rate every 1000 epochs with an initial learning rate of 0.005.

5.2 The CAE

The second module in the DeepMCBM pipeline is a CAE, conditioned on the unwarped panoramic central moments. In our experiments, we conditioned the CAE on the first two such moments, but more generally, any number of moments can be used as well. Theoretically, using more moments implies that the distribution of the pixel stack is better captured.

Our CAE was based on the (unconditional) Autoencoder from [2]. Using a similar architecture, we used 4 channels in each hidden convolutional layer and 2 channels for the last convolutional layer. In between the convolutional encoder and decoder we used dense layers of size 256×4 and 4×256 (thus, the code size was 4) with a ReLU [1] activation function. All convolutional layers used a 5×5 kernel with stride size 2.

Our conditioning is done in both the encoder and decoder parts. As the conditioning (the unwarped central moments) is in the same spatial dimensions as the input image, conditioning on the encoder size was done by a simple concatenation, along the channel dimension, of the input and the unwarped moments. This results in $(N_{\text{moments}}+1)\cdot C$ input channels for the encoder where N_{moments} is the number of the moments used $(N_{\text{moments}}=2$ in our experiments) and C is the number of channels of an input image (C=3) in the RGB case). To condition the decoder, we concatenate the unwaped moments to the (classic) decoder output and then use a small Convolutional Neural Net (CNN) to integrate the decoder's output and the conditioning. This last CNN is composed of 3 convolutional layers with 50 channels for the hidden layers and C channels for the output layer. All three layers use a 3×3 kernel, and a ReLU [1] activation function.

Using the robust reconstruction loss mentioned in the paper, we trained the CAE for 2000 epochs with a learning rate scheduler and a step size of 500 epochs and initial learning rate of 0.001.

6 Using the Matrix Exponential within an STN

Let A be a 3×3 real-valued matrix such that its last row is all zeros:

$$\mathbf{A} = \begin{bmatrix} A_1 & A_2 & A_3 \\ A_4 & A_5 & A_6 \\ 0 & 0 & 0 \end{bmatrix} . \tag{5}$$

Then, a known result (which is also widely-used in computer vision; see, e.g., [8,9]) from the theory of matrix groups is that

$$T \triangleq \exp(A) \tag{6}$$

(i.e., the matrix exponential of A) has the following form,

$$T = \begin{bmatrix} T_1 & T_2 & T_3 \\ T_4 & T_5 & T_6 \\ 0 & 0 & 1 \end{bmatrix}$$
 (7)

where, in addition, we have that $\det T > 0$ (in particular, T is invertible). Consequently, the matrix exponential provides a mapping from the unconstrained linear space \mathbb{R}^6 into the Affine Group (namely, the space of invertible affine transformations – in this case from \mathbb{R}^2 to \mathbb{R}^2). Taken together with the fact that the matrix exponential is a differentiable function, this means that an STN [7] can be easily set to produce only invertible transformations [11,4].

References

- 1. Agarap, A.F.: Deep learning using rectified linear units (relu). arXiv preprint arXiv:1803.08375 (2018) 6, 7
- Ballé, J., Laparra, V., Simoncelli, E.P.: End-to-end optimized image compression. In: ICLR (2017) 7
- 3. Black, M.J., Rangarajan, A.: On the unification of line processes, outlier rejection, and robust statistics with applications in early vision. IJCV (1996) 5
- 4. Chelly, I., Winter, V., Litvak, D., Rosen, D., Freifeld, O.: JA-POLS: a moving-camera background model via joint alignment and partially-overlapping local subspaces. In: CVPR (2020) 1, 4, 5, 7
- Geman, S., McClure, D.E.: Statistical methods for tomographic image reconstruction. In: BISI (1987) 5
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
- Jaderberg, M., Simonyan, K., Zisserman, A., et al.: Spatial transformer networks. In: NIPS (2015) 7
- Lin, D., Grimson, E., Fisher III, J.: Learning visual flows: A Lie algebraic approach. In: CVPR (2009) 7
- Lin, D., Grimson, E., Fisher III, J.: Modeling and estimating persistent motion with geometric flows. In: CVPR (2010) 7
- Perazzi, F., Pont-Tuset, J., McWilliams, B., Van Gool, L., Gross, M., Sorkine-Hornung, A.: A benchmark dataset and evaluation methodology for video object segmentation. In: CVPR (2016) 4
- Skafte Detlefsen, N., Freifeld, O., Hauberg, S.: Deep diffeomorphic transformer networks. In: CVPR (2018) 7