

Space-Time Video Completion

Yonatan Wexler, Eli Shechtman, Michal Irani

Dept. of Computer Science and Applied Mathematics

The Weizmann Institute of Science, Rehovot, 76100 Israel

Abstract

This paper presents a new framework for completion of missing information based on local structures. It poses the task of completion as a global optimization problem with a well-defined objective function and derives a new algorithm to optimize it. Missing values are constrained to form coherent structures with respect to reference examples.

We apply this method to space-time completion of large space-time “holes” in video sequences of complex dynamic scenes. The missing portions are filled in by sampling spatio-temporal patches from the available parts of the video, while enforcing global spatio-temporal consistency between *all* patches in and around the hole. The consistent completion of static scene parts simultaneously with dynamic behaviors leads to realistic looking video sequences and images.

Space-time video completion is useful for a variety of tasks, including, but not limited to: (i) Sophisticated video removal (of undesired static or dynamic objects) by completing the appropriate static or dynamic background information, (ii) Correction of missing/corrupted video frames in old movies, (iii) Modifying a visual story by replacing unwanted elements, (iv) Creation of video textures by extending smaller ones, (v) Creation of complete field-of-view stabilized video, and (vi) As images are one-frame videos, we apply the method to this special case as well.

I. INTRODUCTION

We present a method for *space-time completion* of large space-time “holes” in video sequences of complex dynamic scenes. We follow the spirit of [14] and use non-parametric sampling, while extending it to handle static and dynamic information simultaneously. The missing video portions are filled-in by sampling spatio-temporal patches from other video portions, while enforcing global spatio-temporal consistency between all patches in and around the hole. Global consistency is obtained by posing the problem of video completion/synthesis as a global optimization problem with a *well-defined objective function* and solving it appropriately. The objective function states that the resulting completion should satisfy the following two constraints: (i) Every *local* space-time patch of the video sequence should be similar to some local space-time patch in the remaining parts of the video sequence (the “input data-set”), while (ii) *globally* all these patches must be consistent with each other, both spatially and temporally.

Solving the above optimization problem is not a simple task, especially due to the large dimensionality of video data. However, we exploit the spatio-temporal relations and redundancies to speed and constrain the optimization process in order to obtain realistic-looking video sequences

with complex scene dynamics at reasonable computation times.

Figure 1 shows an example of the task at hand. Given the input video (Fig. 1.a), a space-time hole is specified in the sequence (Fig. 1.b). The algorithm is requested to complete this hole using information from the remainder of the sequence. We assume that the hole is provided to the algorithm. While in all examples here it was marked manually, it can also be the outcome of some segmentation algorithm. The resulting completion and the output sequence are shown in Figs. 1(c).

The goal of this work is close to a few well studied domains. *Texture Synthesis* (e.g. [2], [14], [28]) extends and fills regular fronto-parallel image textures. This is similar to *Image Completion* (e.g., [9], [12]) which aims at filling in large missing image portions. While impressive results have been achieved recently in some very challenging cases (e.g., see [12]), the goal and the proposed algorithms have so far been defined only in a heuristic way. Global inconsistencies often result from independent local decisions taken at independent image positions. For this reason, these algorithms use large image patches in order increase the chances of correct output. The two drawbacks of this approach are that elaborate methods for combining the large patches are needed for hiding inconsistencies [12], [13], [21], and that the dataset needs to be artificially enlarged by including various skewed and scaled replicas that might be needed for completion. When compared to the pioneering work of [14] and its derivatives, the algorithm presented here can be viewed as stating an objective function explicitly. From this angle, the algorithm of [14] makes a greedy decision in each pixel based on the currently available pixels around it. It only uses a directional neighborhood around each pixel. A greedy approach requires the correct decision to be made at every step. Hence, the chances for errors increase rapidly as the gap grows. The work of [9] showed that this may be alleviated by prioritizing the completion order using local image structure. In challenging cases, containing complex scenes and large gaps, the local neighborhood does not hold enough information for a globally correct solution. This is more pronounced in video. Due to motion aliasing, there is little chance that an exact match will be found.

The framework presented here requires that the *whole* neighborhood around each pixel is considered, not just a causal subset of it. Moreover, it considers *all* windows containing each pixel simultaneously, thus effectively using an even larger neighborhood.

Image Inpainting (e.g., [5], [6], [22]) was defined in a principled way as an edge continuation

1. Sample frames from the original sequence:



2. Zoomed in view around the space-time hole:

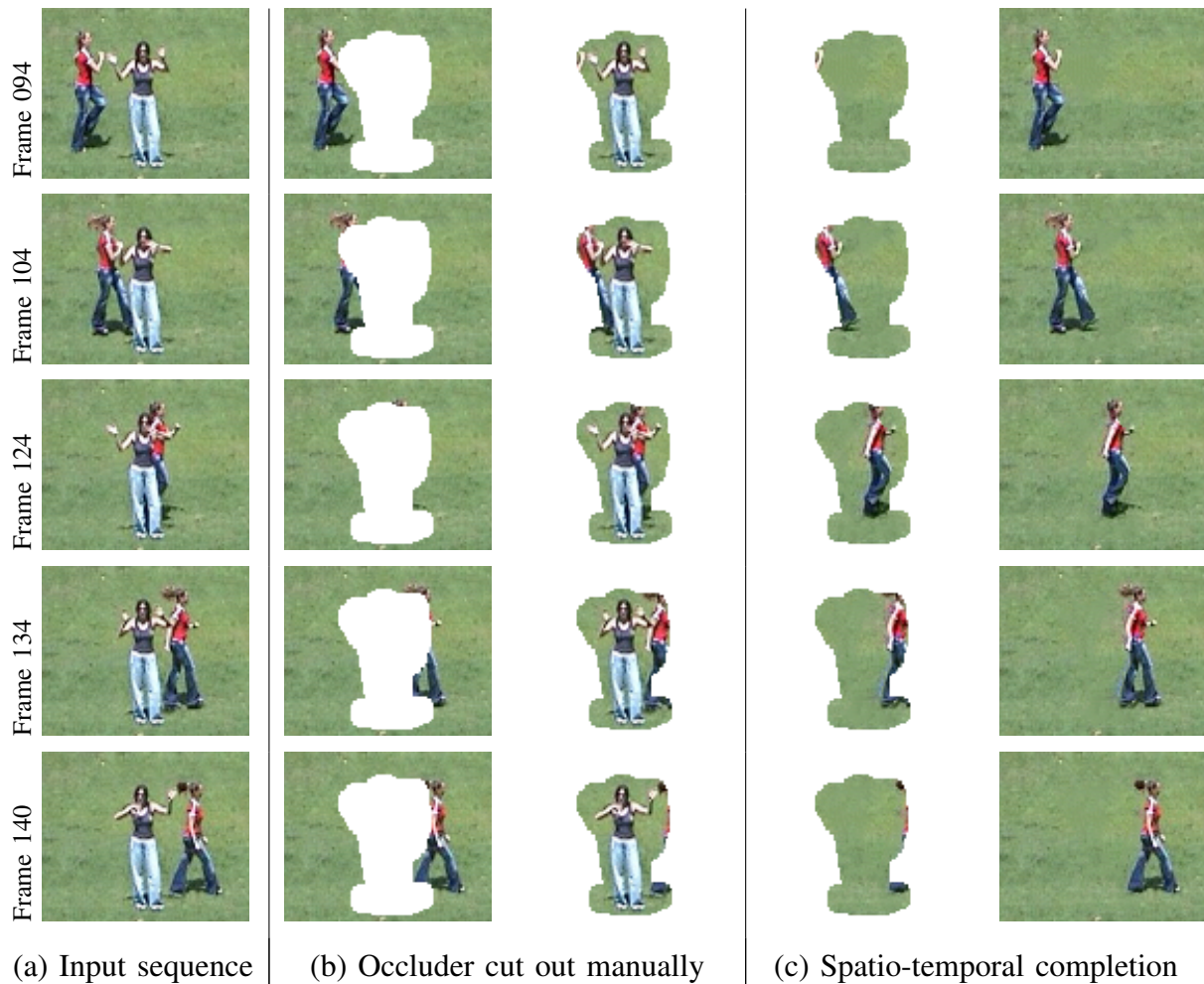


Fig. 1

(1) FEW FRAMES OUT OF A VIDEO SEQUENCE SHOWING ONE PERSON STANDING AND WAVING HER HANDS WHILE THE OTHER PERSON IS HOPPING BEHIND HER. THE VIDEO SIZE IS $100 \times 300 \times 240$, WITH 97,520 MISSING PIXELS. (2) A ZOOMED-IN VIEW ON A PORTION OF THE VIDEO AROUND THE SPACE-TIME HOLE BEFORE AND AFTER THE COMPLETION. NOTE THAT THE RECOVERED ARMS OF THE HOPPING PERSON ARE AT SLIGHTLY DIFFERENT ORIENTATIONS THAN THE REMOVED ONES. AS THIS PARTICULAR INSTANCE DOES NOT EXIST ANYWHERE ELSE IN THE SEQUENCE, A SIMILAR ONE FROM A DIFFERENT TIME INSTANCE WAS CHOSEN TO PROVIDE AN EQUALLY LIKELY COMPLETION.

SEE VIDEO IN: www.wisdom.weizmann.ac.il/~vision/VideoCompletion.html

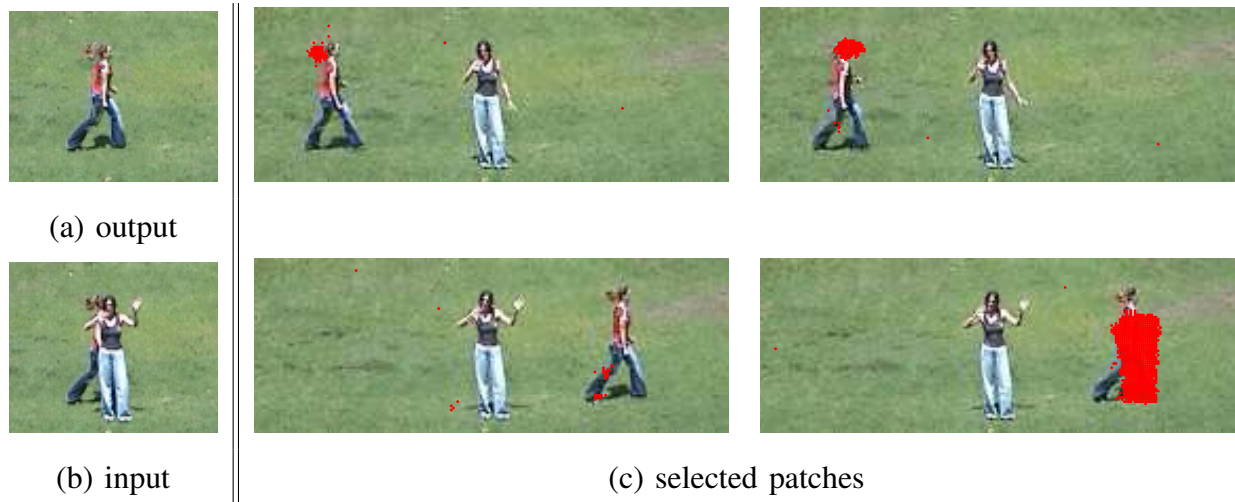


Fig. 2

SOURCES OF INFORMATION. OUTPUT FRAME 114, SHOWN IN (A), IS A COMBINATION OF THE PATCHES MARKED HERE IN RED OVER THE INPUT SEQUENCE (C). IT IS NOTICEABLE THAT LARGE CONTINUOUS REGIONS HAVE BEEN AUTOMATICALLY PICKED WHENEVER POSSIBLE. NOTE THAT THE HAIR, HEAD AND BODY WERE TAKEN FROM DIFFERENT FRAMES. THE ORIGINAL FRAME IS SHOWN IN (B).

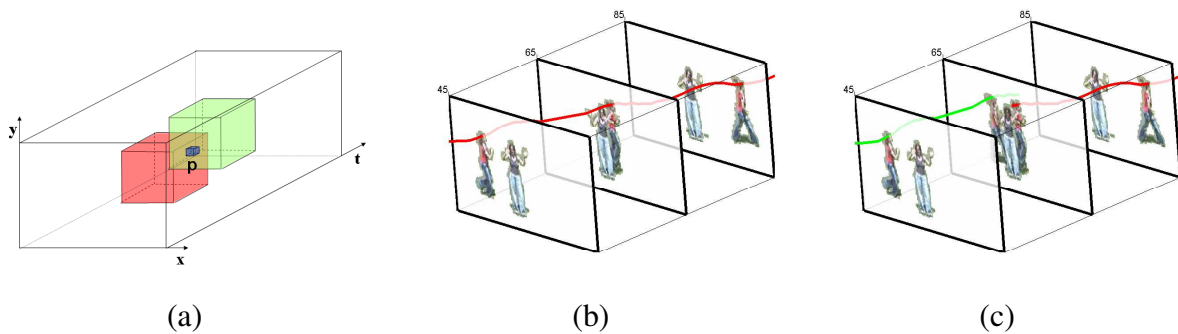


Fig. 3

LOCAL AND GLOBAL SPACE-TIME CONSISTENCY. (A) ENFORCEMENT OF THE GLOBAL OBJECTIVE FUNCTION OF EQ. (1) REQUIRES COHERENCE OF ALL SPACE-TIME PATCHES CONTAINING THE POINT p . SUCH COHERENCE LEADS TO A GLOBALLY CORRECT SOLUTION (B). TRUE TRAJECTORY OF THE MOVING OBJECT IS MARKED IN RED AND IS CORRECTLY RECOVERED. WHEN ONLY LOCAL CONSTRAINTS ARE USED AND NO GLOBAL CONSISTENCY ENFORCED, THE RESULTING COMPLETION LEADS TO INCONSISTENCIES (C). BACKGROUND FIGURE IS WRONGLY RECOVERED TWICE WITH WRONG MOTION TRAJECTORIES.

process, but is restricted to small (narrow) missing image portions in highly structured image data. These approaches have been restricted to completion of *spatial information* alone in images. Even when applied to video sequences (as in [4]), the completion was still performed spatially. The temporal component of video has mostly been ignored. The basic assumption of Image Inpainting, that edges should be interpolated in some smooth way, does not naturally extend to time. Temporal aliasing is typically much stronger than spatial aliasing in video sequences of dynamic scenes. Often a pixel may contain background information in one frame and foreground information in the next frame, resulting in very non-smooth temporal changes. These violate the underlying assumptions of Inpainting.

In [19] a method has been proposed for employing spatio-temporal information to correct scratches and noise in poor-quality video sequences. This approach relies on optical-flow estimation and propagation into the missing parts followed by reconstruction of the color data. The method was extended to removal of large objects in [20] under the assumption of planar rigid layers and small camera motions.

View Synthesis deals with computing the appearance of a scene from a new view point, given several images. A similar objective was already used in [15] for successfully resolving ambiguities which are otherwise inherent to the geometric problem of new-view synthesis from multiple camera views. The objective function of [15] was defined on 2D images. Their local distance between 2D patches was based on SSD of color information and included geometric constraints. The algorithm there did not take into account the dependencies between neighboring pixels as only the central point was updated in each step.

Recently there have been a few notable approaches to video completion. From the algorithmic perspective, the most similar work to ours is [7] which learns a mapping from the input video into a smaller volume, aptly called “epitome”. These are then used for various tasks including completion. While our work has a similar formulation, there are major differences. We seek some cover of the missing data by the available information. Some regions may contribute more than once while some may not contribute at all. In contrast, the epitome contains a proportional representation of all the data. One implication of this is that the windows will be averaged and so will lose some fidelity. The recent work of [24] uses estimated optical flow to separate foreground and background layers. Each is then filled incrementally using the priority-based ordering idea similar to [9].

Lastly, the work of [18] takes an object-based approach where large portions of the video are tracked, their cycles are analyzed and can then be inserted into the video. This allows warping of the object so it fits better, but requires a complete appearance of the object to be identified and so is not applicable to more complex dynamics, such as articulated motion or stochastic textures.

A closely related area of research which regards temporal information explicitly is that of *dynamic texture synthesis* in videos (e.g. [3], [11]). Dynamic textures are usually characterized by an unstructured stochastic process. They model and synthesize smoke, water, fire, etc. but cannot model nor synthesize structured dynamic objects, such as behaving people. We demonstrate the use of our method for synthesizing large video textures from small ones in section VI. While [26] has been able to synthesize/complete video frames of structured dynamic scenes, it assumes that the “missing” frames already appear in their entirety elsewhere in the input video, and therefore needed only to identify the correct permutation of frames. An extension of that paper [25] manually composed smaller “video sprites” to a new sequence.

The approach presented here can automatically handle the completion and synthesis of both structured dynamic objects as well as unstructured dynamic objects under a single framework. It can complete frames (or portions of them) that never existed in the dataset. Such frames are constructed from various space-time patches, which are automatically selected from different parts of the video sequence, all put together consistently. The use of a global objective function removes the pitfalls of local inconsistencies and the heuristics of using large patches. As can be seen in the figures and in the attached videos, the method is capable of completing large space-time areas of missing information containing complex structured dynamic scenes, just as it can work on complex images. Moreover, this method provides a unified framework for various types of image and video completion and synthesis tasks, with the appropriate choice of the spatial and temporal extents of the space-time “hole” and of the space-time patches. A preliminary version of this work appeared in CVPR 04’ [29].

The paper is organized as follows: Section II introduces the objective function and Section III describes the algorithm used for optimizing it. Sections IV and V discuss tradeoffs between space and time dimensions in video and how they are unified into one framework. Section VI demonstrates the application of this method to various problems. Finally, Section VII concludes this paper.

II. COMPLETION AS A GLOBAL OPTIMIZATION

Given a video sequence, we wish to complete the missing portions in such a way that it looks just like the available parts. For this we define a global objective function to rank the quality of a completion. Such a function needs to take a completed video and rate its quality with respect to a reference one. Its extremum should denote the best possible solution.

Our basic observation is that in order for a video to look good, it needs to be coherent everywhere. That is, a good completion should resemble the given data locally everywhere.

The constraint on the color of each pixel depends on the joint color assignment of its neighbors. This induces a structure where each pixel depends on the neighboring ones. The correctness of a pixel value depends on whether its local neighborhood forms a coherent structure. Rather than modeling this structure, we follow [14] and use the reference video as a library of video samples that are considered to be coherent.

Given these guidelines, we are now ready to describe our framework in detail.

A. The global objective function

To allow for a uniform treatment of dynamic and static information, we treat video sequences as space-time volumes. We use the following notations. A pixel (x, y) in a frame t will be regarded as a space-time point $p = (x, y, t)$ in the volume. W_p denotes a small, fixed-sized window around the point p both in space and in time. The diameter of the window is given as a parameter. We use indices i and j to denote locations relative to p . For example, W_p is the window centered around p and W_p^i is the i^{th} window containing p and is centered around the i^{th} neighbor of p .

We say that a video sequence \mathcal{S} has *global visual coherence* with some other sequence \mathcal{D} if every local space-time patch in \mathcal{S} can be found somewhere within the sequence \mathcal{D} . In other words, we can cover \mathcal{S} with small windows from \mathcal{D} . Windows in the dataset \mathcal{D} are denoted by V and are indexed by a reference pixel (e.g. V_q and V_q^i).

Let \mathcal{S} and $\mathcal{H} \subseteq \mathcal{S}$ be an input sequence and a ‘‘hole’’ region within it. That is, \mathcal{H} denotes all the missing space-time points within \mathcal{S} . For example, \mathcal{H} can be an undesired object to be erased, a scratch or noise in old corrupt footage, or entire missing frames, etc. We assume that both \mathcal{S} and \mathcal{H} are given.

We wish to complete the missing space-time region \mathcal{H} with some new data \mathcal{H}^* such that the resulting video sequence \mathcal{S}^* will have as much global visual coherence with some reference sequence \mathcal{D} (the dataset). Typically, $\mathcal{D} = \mathcal{S} \setminus \mathcal{H}$, namely - the remaining video portions outside the hole are used to fill in the hole. Therefore, we seek a sequence \mathcal{S}^* which maximizes the following objective function:

$$\text{Coherence}(\mathcal{S}^*|\mathcal{D}) = \prod_{p \in \mathcal{S}^*} \max_{q \in \mathcal{D}} \text{sim}(W_p, V_q) \quad (1)$$

where p, q run over all space-time points in their respective sequences. $\text{sim}(\cdot, \cdot)$ is a local similarity measure between two space-time patches that will be defined shortly (section II-B). The patches need not necessarily be isotropic and can have different sizes in the spatial and temporal dimensions. We typically use $5 \times 5 \times 5$ patches which are large enough to be statistically meaningful but small enough so that effects, such as parallax or small rotations, will not affect them. They are the basic building blocks of the algorithm. The use of windows with temporal extent assumes that the patches are correlated in time, and not only in space. When the camera is static or fully stabilized, this is trivially true. However it may also apply in other cases where the same camera motion appears in the database (as shown in Fig.16).

Figure 3 explains why Eq.(1) induces global coherence. Each space-time point p belongs to other space-time patches of other space-time points in its vicinity. For example, for a $5 \times 5 \times 5$ window, 125 different patches involve p . The red and green boxes in Fig. 3(a) are examples of two such patches. Eq. (1) requires that all 125 patches agree on the value of p and therefore Eq. (1) leads to globally coherent completions such as the one in Fig.3(b). If the global coherence of Eq. (1) is not enforced, and the value of p is determined locally by a single best-matching patch (i.e. using a sequential greedy algorithm as in [9], [14], [24]), then global inconsistencies will occur in a later stage of the recovery. An example of temporal incoherence is shown in Fig. 3(c). A greedy approach requires the correct answer to be found at every step and this is rarely the case as local image structure will often contain ambiguous information.

The above objective function seeks a cover of the missing data using the available one. That is, among the exponentially large number of possible covers, we seek the one that will give us the least amount of total error. The motivation here is to find a solution that is correct everywhere, as any local inconsistency will push the solution away from the minimum.

B. The local space-time similarity measure

At the heart of the algorithm is a well-suited similarity measure between space-time patches. A good measure needs to agree *perceptually* with a human observer. The Sum of Squared Differences (SSD) of color information, that is so widely used for image completion, does not suffice for the desired results in video (regardless of the choice of color space). The main reason for this is that the human eye is very sensitive to motion. Maintaining motion continuity is more important than finding the exact spatial pattern match within an image of the video.

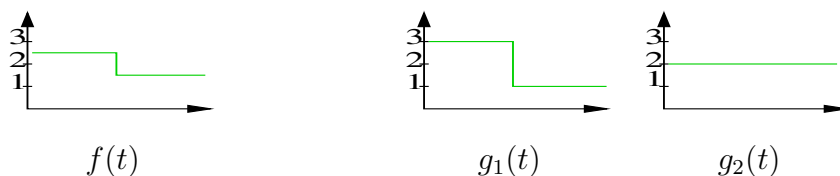


Fig. 4

IMPORTANCE OF MATCHING DERIVATIVES (SEE TEXT).

Figure 4 illustrates (in 1D) that very different temporal behaviors can lead to the same SSD score. The function $f(t)$ has a noticeable temporal change. Yet, its SSD score relative to a similar-looking function $g_1(t)$ is the same as the SSD score of $f(t)$ with a flat function $g_2(t)$: $\int (f - g_1)^2 dt = \int (f - g_2)^2 dt$. However, perceptually, $f(t)$ and $g_1(t)$ are more similar as they both encode a temporal change.

We would like to incorporate this into our algorithm so to create a similarity measure that agrees with perceptual similarity. Therefore we add a measure that is similar to that of normal-flow to obtain a quick and rough approximation of the motion information. Let Y be the sequence containing the grayscale (intensity) information obtained from the color sequence. At each space-time point we compute the spatial and temporal derivatives (Y_x, Y_y, Y_t) . If the motion were only horizontal, then $u = Y_t/Y_x$ would capture the instantaneous motion in the x direction. If the motion were only vertical, then $v = Y_t/Y_y$ would capture the instantaneous motion in the y direction. The fractions factor out most of the dependence between the spatial and the temporal changes (such as frame-rate) while capturing object velocities and directions. These two components are scaled and added to the RGB measurements to obtain a 5-dimensional representation for each space-time point: $(R, G, B, \alpha u, \alpha v)$ where $\alpha = 5$. Note that we do not

compute optical flow. We apply an L_2 -norm (SSD) to this 5-D representation in order to capture space-time similarities for static and dynamic parts simultaneously. Namely, for two space-time windows W_p and V_q we have $d(W_p, V_q) = \sum_{(x,y,t)} \|W_p(x, y, t) - V_q(x, y, t)\|^2$ where for each (x, y, t) within the patch $W_p(x, y, t)$ denotes its 5D measurement vector. The distance is translated to the similarity measure

$$\text{sim}(W_p, V_q) = e^{-\frac{d(W_p, V_q)}{2\sigma^2}} \quad (2)$$

The choice of σ is important as it controls the smoothness of the induced error surface. Rather than using a fixed value, it is chosen to be the 75-percentile of all distances in the current search in all locations. In this way, the majority of locations are taken into account and hence, there is a high probability that the overall error will reduce.

III. THE OPTIMIZATION

The inputs to the optimization are a sequence \mathcal{S} and a ‘‘hole’’ $\mathcal{H} \subset \mathcal{S}$, marking the missing space-time points to be corrected or filled-in. The algorithm seeks an assignment of color values for all the space-time points (pixels) in the hole so to satisfy Eq. (1). While Eq. (1) does not imply any optimization scheme, we note that it will be satisfied if the following two local conditions are met at every space-time point p :

(i) All windows $W_p^1 \dots W_p^k$ containing p appear in the dataset \mathcal{D} :

$$\exists V^i \in \mathcal{D}, W_p^i = V^i$$

(ii) All those $V^1 \dots V^k$ agree on the color value c at location p :

$$c = V^i(p) = V^j(p)$$

This is trivially true since the second condition is a particular case of the first. These conditions imply an iterative algorithm. The iterative step will aim at satisfying these two conditions locally at every point p . Any change in p will affect all the windows that contain it and so the update rule must take all of them into account.

Note that Eq. (1) may have an almost trivial solution, in which the hole \mathcal{H} contains an exact copy of some part in the database. For such a solution, the error inside the hole will be zero (as both conditions above are satisfied) and so the total coherence error will be proportional to the surface area of the space-time boundary. This error might be much smaller than small noise errors spread across the entire volume of the hole. Therefore, we associate an additional

quantity α_p to each point $p \in \mathcal{S}$. Known points $p \in \mathcal{S} \setminus \mathcal{H}$ have fixed high confidence, whereas missing points $p \in \mathcal{H}$ will have lower confidence. This weight is chosen so to ensure that the total error inside the hole is less than that on the hole boundary. This argument needs to hold recursively in every layer in the hole. An approximation to such weighting is to compute the distance transform for every hole pixel, and then use $\alpha_p = \gamma^{-\text{dist}}$. When the hole is roughly spherical choosing $\gamma = 1.3$ gives the desired weighting. This measure bears some similarity to the priority used in [9] except that here it is fixed throughout the process. Another motivation for using such weighting is to speed up convergence by directing the flow of information inwards from the boundary.

A. The iterative step

Let $p \in \mathcal{H}$ be some hole point which we wish to improve. Let $W_p^1 \dots W_p^k$ be all space-time patches containing p . Let $V^1 \dots V^k$ denote the patches in \mathcal{D} that are most similar to $W_p^1 \dots W_p^k$ per Eq. (2). According to condition (i) above, if W_p^i is reliable then $d(W_p^i, V^i) \approx 0$. Therefore $\text{sim}(W_p^i, V^i)$ measures the degree of reliability of the patch W_p^i .

We need to pick a color c at location p so that the coherence at all windows containing it will increase. Each window V^i provides an evidence of possible solution for c and the confidence in this evidence is given by Eq. (2) as $s_p^i = \text{sim}(W_p^i, V^i)$. According to condition (ii) above, the most likely color c at p should minimize the variance of the colors $c^1 \dots c^k$ proposed by $V^1 \dots V^k$ at location p . As mentioned before, these values should be scaled according to the fixed α_p^i to avoid the trivial solution. Thus, the most likely color at p will minimize $\sum_i \omega_p^i (c - c^i)^2$ where $\omega_p^i = \alpha_p^i \cdot s_p^i$. Therefore:

$$c = \frac{\sum_i \omega_p^i c^i}{\sum_i \omega_p^i} \quad (3)$$

This c is assigned to be the color of p at the end of the current iteration. Such an update rule minimizes the local error around each space-time point $p \in \mathcal{H}$ while maintaining consistency in all directions, leading to a global consistency.

One drawback of the update rule in Eq. (3) is its sensitivity to outliers. It is enough that a few neighbors suggest the wrong color to bias the mean color c and thus prevent or delay the convergence of the algorithm. In order to avoid such effects, we treat the k possible assignments as evidence. These evidences give discrete samples in the continuous space of possible color

- 1: **Input:** video \mathcal{S} , hole $\mathcal{H} \subset \mathcal{S}$, database \mathcal{D} .
- 2: Compute space-time pyramids $\mathcal{S}_l, \mathcal{H}_l, \mathcal{D}_l$ $l = 1..L$.
- 3: $t \leftarrow 0, \mathcal{S}^t \leftarrow \mathcal{S}$
- 4: **for** pyramid level l , from top to bottom **do**
- 5: **repeat**
- 6: **for all** $p \in \mathcal{H}_l$ **do**
- 7: Let $\{W_p^i\}_{i=1}^k$ be all windows s.t. $p \in W_p^i$
- 8: Find $\{V^i\} \subseteq \mathcal{D}_l$ maximizing Eq. (2)
- 9: Let $c^i \in V^i$ be the appropriate colors.
- 10: Set $\omega_p^i = \alpha_p^i \cdot \text{sim}(W_p^i, V^i)$.
- 11: $\mathcal{S}^{t+1}(p) \leftarrow ML(c^i, \omega_p^i)$ using Eq. (4)
- 12: **end for**
- 13: $t \leftarrow t + 1$
- 14: **until** convergence
- 15: Propagate solution to the next level
- 16: **end for**
- 17: **Output:** \mathcal{S}^t

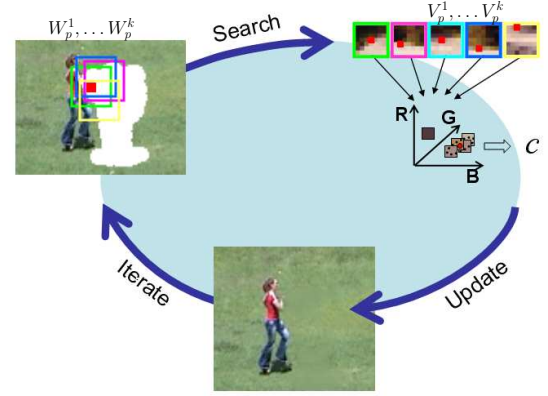


Fig. 5

VIDEO COMPLETION ALGORITHM

assignments. The reliability of each evidence is proportional to ω_p^i and we seek the Maximum Likelihood (ML) in this space. This is computed using the Mean-Shift algorithm [8] with a variable bandwidth (window size). The Mean-Shift algorithm finds the density modes of the distribution (which is related to Parzen windows density estimation). It is used here to extract the dominant mode of the density. The bandwidth is defined with respect to the standard deviation σ of the colors c_1, \dots, c_k at each point. It is typically set to be 3σ at the beginning and is reduced gradually to 0.2σ . The highest mode M is picked and so the update rule in Table 5 is:

$$c = \frac{\sum_{i \in M} \omega_p^i c^i}{\sum_{i \in M} \omega_p^i} \quad (4)$$

This update rule produces a robust estimate for the pixel in the presence of noise. While making the algorithm slightly more complex, this step has further advantages. The bandwidth parameter controls the degree of smoothness of the error surface. When it is large, it reduces to simple

weighted averaging (as in Eq. 3), hence allowing rapid convergence. When it is small, it induces a weighted majority vote, avoiding blurring in the resulting output. The use of varying σ value has significantly improved the convergence of the algorithm. When compared with the original work in [29], the results shown here achieve better convergence. This can be seen from the improved amount of details in areas which were previously smoothed unnecessarily. The ability to rely on outlier rejection also allows us to use approximate nearest neighbors, as discussed in Sec.III-C

B. Multi-scale solution

To further enforce global consistency and to speed up convergence, we perform the iterative process in multiple scales using spatio-temporal pyramids. Each pyramid level contains half the resolution in the spatial and in the temporal dimensions. The optimization starts at the coarsest pyramid level and the solution is propagated to finer levels for further refinement. Figure 6 shows a typical multiscale V-cycle performed by the algorithm. It is worth to mention that as each level contains $1/8^{\text{th}}$ of the pixels, both in the hole \mathcal{H} and in the database \mathcal{D} , the computational cost of using a pyramid is almost negligible ($8/7$ of the work). In hard examples, it is sometimes necessary to repeat several such cycles, gradually reducing the pyramid height. This is inspired by multi-grid methods [27].

The propagation of the solution from a coarse level to the one above it is done as follows. Let p_{\uparrow} be a location in the finer level and let p_{\downarrow} be its corresponding location in the coarser level. As before, let $W_{p_{\downarrow}}^i$ be the windows around p_{\downarrow} and let V_{\downarrow}^i be the matching windows in the database. We propagate the locations of V_{\downarrow}^i onto the finer level to get V_{\uparrow}^i . Some of these (about $\frac{k}{8}$, half in each dimension) will overlap p_{\uparrow} and these will be used for the maximum-likelihood step, just as before (except that here there are less windows). This method is better than plain interpolation as the initial guess for the next level will preserve high spatio-temporal frequencies and will not be blurred unnecessarily.

C. Algorithm Complexity

We now discuss the computational complexity of the suggested method. Assume we have $N = |\mathcal{H}|$ pixels in the hole and that there are K windows in the database. The algorithm has the following structure.

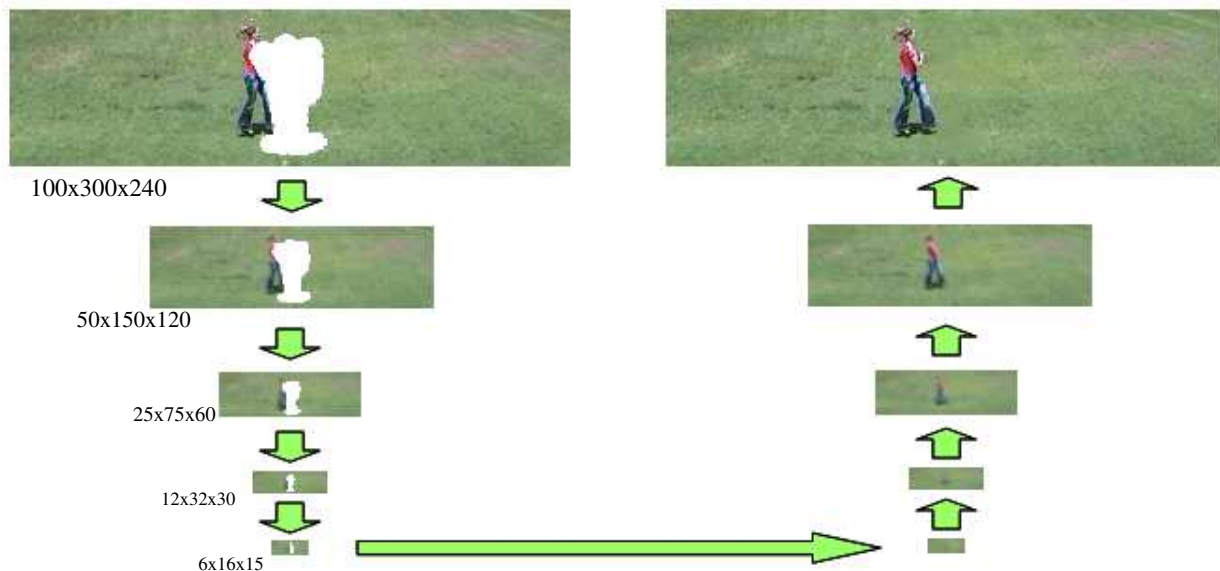


Fig. 6

MULTISCALE SOLUTION. THE ALGORITHM STARTS BY SHRINKING THE INPUT VIDEO (BOTH IN SPACE AND IN TIME). THE COMPLETION STARTS AT THE COARSEST LEVEL, ITERATING EACH LEVEL SEVERAL TIMES AND PROPAGATING THE RESULT UPWARDS.

First a spatio-temporal pyramid is constructed and the algorithm iterates a few times in each level (typically five to ten iterations in each level). Level l contains roughly $N/8^l$ hole pixels, and $K/8^l$ database windows.

Each iteration of the algorithm has two stages: Searching the database and per-pixel Maximum Likelihood.

The first stage performs a nearest-neighbor search once for each window overlapping the space-time hole \mathcal{H} . This is the same computational cost as all the derivatives from Efros' work [14]. The search time depends on K and the nearest-neighbor search method. Since each patch is searched independently, this step can be parallelized trivially. While brute-force would take $O(K/8^l \cdot N/8^l)$, we use the method of [1] so the search time is logarithmic in K . We typically obtain a speedup of two orders of magnitude over brute force search. This approach is very suitable for our method for three reasons. First, it is much faster. Second, we always search for windows of the same size. Third, our robust maximum likelihood estimation allows us to deal

with wrong results that may be returned by this approximate algorithm.

The second stage is the per pixel maximum likelihood computation. We do this using the Mean-Shift algorithm [8]. The input is a set of 125 RGB triplets along with their weights. A trivial implementation is quadratic in this small number and so is very fast.

The running time depends on the above factors. For small problems, such as image completion, the running time is about a minute. For the ‘‘Umbrella’’ sequence in Fig. 7 of size $120 \times 340 \times 100$, with 422,000 missing pixels each iteration in the top level takes roughly one hour on a 2.6Ghz Pentium computer. Roughly 95% of this time is used for the nearest neighbor search. This suggests that pruning the database (as in [24]) or simplifying the search (as in [2]) would give a significant speedup.

D. Relation to statistical estimation

The above algorithm can be viewed in terms of a statistical estimation framework. The global visual coherence in Eq.(1) can be derived as a Likelihood function via a graphical model, and our iterative optimization process can be seen as an approximation of the EM method to find the maximum likelihood estimate. According to this model the parameters are the colors of the missing points in the space-time hole and pixels in the boundary around the hole are the observed variables. The space-time patches are the hidden variables and are drawn from the dataset (patches outside the hole in our case).

We show in the Appendix how the likelihood function is derived from this model, and that the maximum likelihood solution is the best visually coherent completion, as defined in Eq. (1). We also show that our optimization algorithm fits the EM method [10] for maximizing the Likelihood function. Under some simplifying assumptions the **E** step is equivalent to the nearest neighbor match of a patch from the dataset to the current estimated corresponding ‘‘hole’’ patch. The **M** step is equivalent to the update rule of Eq. 3.

The above procedure also bears some similarity to the Belief Propagation (BP) approach to completion such as in [22]. As BP communicates a PDF (probability density function), it is practically limited to modeling no more than three neighboring connections at once. There, a standard grid-based graphical model is used (as is for example in [16]), in which each pixel is connected to its immediate neighbors. Our graphical model has a *completely different* structure. Unlike BP our model does not have links between nodes (pixels). Rather, the patch structure

induces an implicit connection between nodes. It takes into account not only a few immediate neighbors but *all* of them (e.g. 125). This allows us to deal with more complex visual structures. In addition, we assume that the point values are parameters and not random variables. Instead of modeling the PDF, we use the local structure to estimate its likelihood and derive an update step using the nearest neighbor in the dataset.

When seen as a variant of belief propagation, one can think of this method as propagating only one evidence. In the context of this work, this is reasonable as the space of all patches is very high dimensional and is sparsely populated. Typically, the dataset occupies “only” few million samples out of $(3 * 255)^{125}$ possible combinations. The typical distance between the points is very high and so only a few are relevant, so passing a full PDF is not advantageous.

IV. SPACE-TIME VISUAL TRADEOFFS

The spatial and temporal dimensions are very different in nature, yet are inter-related. This introduces visual tradeoffs between space and time that are beneficial to our space-time completion process. On one hand, these relations are exploited to narrow down the search space and to speed up the completion process. On the other hand, they often entail different treatments of the spatial and temporal dimensions in the completion process. Some of these issues have been mentioned in previous sections in different contexts, and are therefore only briefly mentioned here. Other issues are discussed here in more length.

Temporal vs. spatial aliasing: Typically, there is more temporal aliasing than spatial aliasing in video sequences of dynamic scenes. This is mainly due to the different nature of blur functions that precede the sampling process (digitization) in the spatial and in the temporal dimensions: The spatial blur induced by the video camera (a Gaussian whose extent is several pixels) is a much better low-pass filter than the temporal blur induced by the exposure time of the camera (a Rectangular blur function whose extent is less than a single frame-gap in time). This leads to a number of observations:

1. Extending the family of Inpainting methods to include the temporal dimension may be able to handle completion of (narrow) missing video portions that undergo slow motions, but there is a high likelihood that it will not be able to handle fast motions or even simple everyday human motions (such as walking, running, etc). This is because Inpainting relies on edge continuity,

which will be hampered by strong temporal aliasing.

Space-time completion, on the other hand, does not rely on smoothness of information within patches, and can therefore handle aliased data as well.

2. Because temporal aliasing is shorter than spatial aliasing, our multi-scale treatment is not identical in space and in time. In particular, after applying the video completion algorithm of Sec. III, residual spatial (appearance) errors may still appear in fast recovered moving objects. To correct those effects, an additional refinement step of space-time completion is added, but this time only the spatial scales vary (using a spatial pyramid), while the temporal scale is kept at the original temporal resolution. The completion process, however, is still space-time. This allows for completion using patches which have a large spatial extent, to correct the spatial information, while maintaining a minimal temporal extent so that temporal coherence is preserved without being affected too much by the temporal aliasing.

The local patch size: In our space-time completion process we typically use $5 \times 5 \times 5$ patches. Such a patch size provides $5^3 = 125$ measurements per patch. This usually provides sufficient statistical information to make reliable inferences based on this patch. To obtain a similar number of measurements for reliable inference in the case of 2D image completion, we would need to use patches of size 11×11 . Such patches, however, are not small, and are therefore more sensitive to geometric distortions (effects of 3D parallax, change in scale and orientation) due to different viewing directions between the camera and the imaged objects. This restricts the applicability of image-based completion, or else requires the use of patches at different sizes and orientations [12], which increases the complexity of the search space combinatorially.

One may claim that due to the new added dimension (time) there is a need to select patches with a larger number of samples, to reflect the increase in data complexity. This, however, is not the case, due to the large degree of spatio-temporal redundancy in video data. The data complexity indeed increases slightly, but this increase is in no way proportional to the increase in the amounts of data.

The added temporal dimension therefore provides greater flexibility. The locality of the $5 \times 5 \times 5$ patches both in space and in time makes them relatively insensitive to small variations in scaling, orientation, or viewing direction, and therefore applicable to a richer set of scenes (richer both in the spatial sense and in the temporal sense).

Interplay between time and space: Often the lack of spatial information can be compensated for by the existence of temporal information, and vice versa. To show the importance of combining the two cues of information, we compare the results of spatial completion alone to those of space-time completion. The top row of Fig. 9 displays the resulting completed frames of Fig. 1 using space-time completion. The bottom row of Fig. 9 shows the results obtained by filling-in the same missing regions, but this time using only image (spatial) completion. In order to provide the 2D image completion with the best possible conditions, the image completion process was allowed to choose the spatial image patches from *any* of the frames in the input sequence. It is clear from the comparison that image completion failed to recover the dynamic information. Moreover, it failed to complete the hopping woman in any reasonable way, regardless of the temporal coherence.

Furthermore, due to the large spatio-temporal redundancy in video data, the added temporal component provides additional flexibility in the completion process. When the missing space-time region (the “hole”) is spatially large and temporally small, then the temporal information will provide most of the constraints in the completion process. In such cases image completion will not work, especially if the missing information is dynamic. Similarly, if the hole is temporally large, but spatially small, then spatial information will provide most of the constraints in the completion, whereas pure temporal completion/synthesis will fail. Our approach provides a unified treatment of all these cases, without the need to commit to a spatial or a temporal treatment in advance.

V. UNIFIED APPROACH TO COMPLETION

The approach presented in this paper provides a unified framework for various types of image and video completion/synthesis tasks. With the appropriate choice of the spatial and temporal extents of the space-time “hole” \mathcal{H} and of the space-time patches W_p , our method reduces to any of the following special cases:

1. When the space-time patches W_p of Eq. (1) have only a spatial extent (i.e., their temporal extent is set to 1), then our method becomes the classical *spatial* image completion and synthesis. However, because our completion process employs a global objective function (Eq. 1), global consistency is obtained that is otherwise lacking when not enforced. A comparison of our method to other image completion/synthesis methods is shown in Figs 12, 13 and 14. (We could not

check the performance of [12] on these examples. We have, however, applied our method to the examples shown in [12], and obtained comparably good results.)

2. When the spatial extent of the space-time patches W_p of Eq. (1) is set to be the entire image, then our method reduces to *temporal* completion of missing frames or synthesis of new frames using existing frames to fill in temporal gaps (similar to the problem posed by [26]).

3. If, on the other hand, the spatial extent of the space-time “hole” \mathcal{H} is set to be the entire frame (but the patches W_p of Eq. (1) remain small), then our method still reduces to *temporal* completion of missing video frames (or synthesis of new frames), but this time, unlike [26], the completed frames may have never appeared in their entirety anywhere in the input sequence. Such an example is shown in Fig. 8, where three frames were dropped from the video sequence of a man walking on the beach. The completed frames were synthesized from bits of information gathered from different portions of the remaining video sequence. Waves, body, legs, arms, etc., were automatically selected from different space-time locations in the sequence so that they all match coherently (both in space and in time) to each other as well as to the surrounding frames.

VI. APPLICATIONS

Space-time video completion is useful for a variety of tasks in video post production and video restoration. A few example applications are listed below. In all the examples below we crop the video so to contain only relevant portions. The size of the original videos for most of them is half PAL video resolution (360×288).

1. Sophisticated video removal: Video sequences often contain undesired objects (static or dynamic), which were either not noticed or else were unavoidable at the time of recording. When a moving object reveals all portions of the background at different time instances, then it can be easily removed from the video data, while the background information can be correctly recovered using simple techniques (e.g., [17]). Our approach can handle the more complicated case, when portions of the background scene are never revealed, and these occluded portions may further change dynamically. Such examples are shown in Figs. 1, 7, and 10. Note that both in 7, and 10, all the completed information is dynamic and so reliance on background subtraction or segmentation is not likely to succeed here.

2. Restoration of old movies: Old video footage is often very noisy and visually corrupted. Entire frames or portions of frames may be missing, and severe noise may appear in other video

portions. These kinds of problems can be handled by our method. Such an example is shown in Fig 11.

(a) Input sequence



(b) Umbrella removed



(c) Output



(d) Full-size frames (input and output):



Fig. 7

REMOVAL OF A BEACH UMBRELLA. THE VIDEO SIZE IS $120 \times 340 \times 100$, WITH 422,000 MISSING PIXELS. SEE VIDEO IN:

www.wisdom.weizmann.ac.il/~vision/VideoCompletion.html

Input video with seven missing frames:



Output video with completed frames:

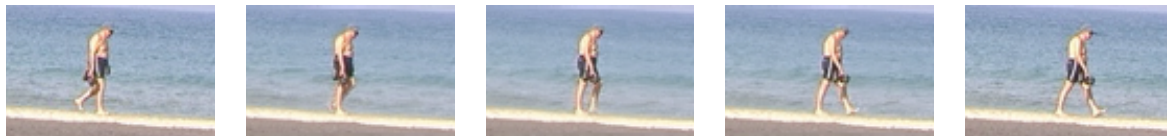


Fig. 8

COMPLETION OF MISSING FRAMES. THE VIDEO SIZE IS $63 \times 131 \times 42$ AND HAS 57,771 MISSING PIXELS SEVEN FRAMES.

Video Completion:

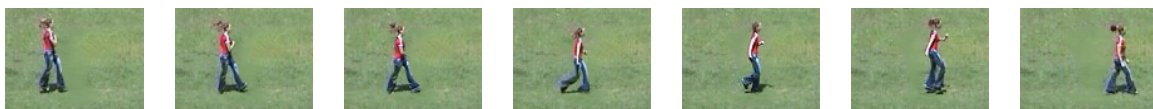


Image Completion:



Fig. 9

IMAGE COMPLETION VERSUS VIDEO COMPLETION

Input video:



One lady cut out:



Output video with completed frames:



Fig. 10

REMOVAL OF A RUNNING PERSON. THE VIDEO SIZE IS $80 \times 170 \times 88$ WITH 74,625 MISSING PIXELS.



Fig. 11

RESTORATION OF A CORRUPTED OLD CHARLIE CHAPLIN MOVIE. THE VIDEO IS OF SIZE $135 \times 180 \times 96$ WITH 3,522 MISSING PIXELS.

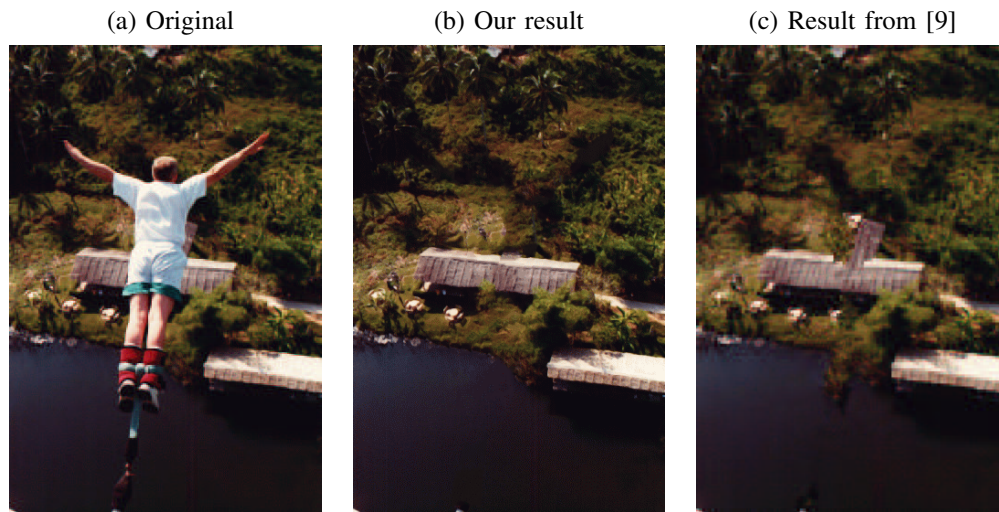


Fig. 12

IMAGE COMPLETION EXAMPLE.

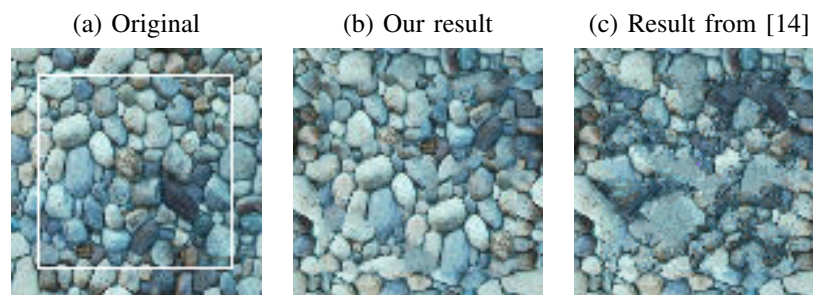


Fig. 13

TEXTURE SYNTHESIS EXAMPLE. THE ORIGINAL TEXTURE (A) WAS ROTATED 90° AND ITS CENTRAL REGION (MARKED BY WHITE SQUARE) WAS FILLED USING THE ORIGINAL INPUT IMAGE. (B) OUR RESULT. (C) BEST RESULTS OF (OUR IMPLEMENTATION OF) [14] WERE ACHIEVED WITH 9×9 WINDOWS (THE ONLY USER DEFINED PARAMETER). ALTHOUGH THE WINDOW SIZE IS LARGE, THE GLOBAL STRUCTURE IS NOT MAINTAINED.

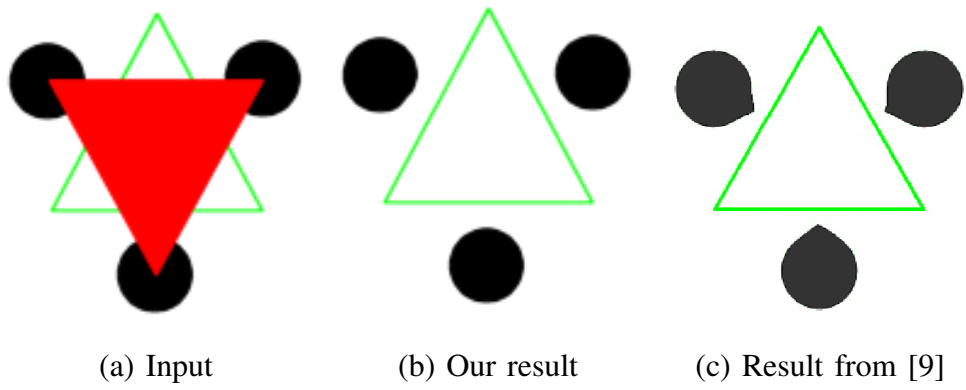


Fig. 14

THE KANITZSA TRIANGLE EXAMPLE IS TAKEN FROM [9] AND COMPARED WITH THE ALGORITHM HERE. THE PREFERENCE OF [9] TO COMPLETE STRAIGHT LINES CREATES THE BLACK CORNERS IN (C) WHICH DO NOT APPEAR IN (B).

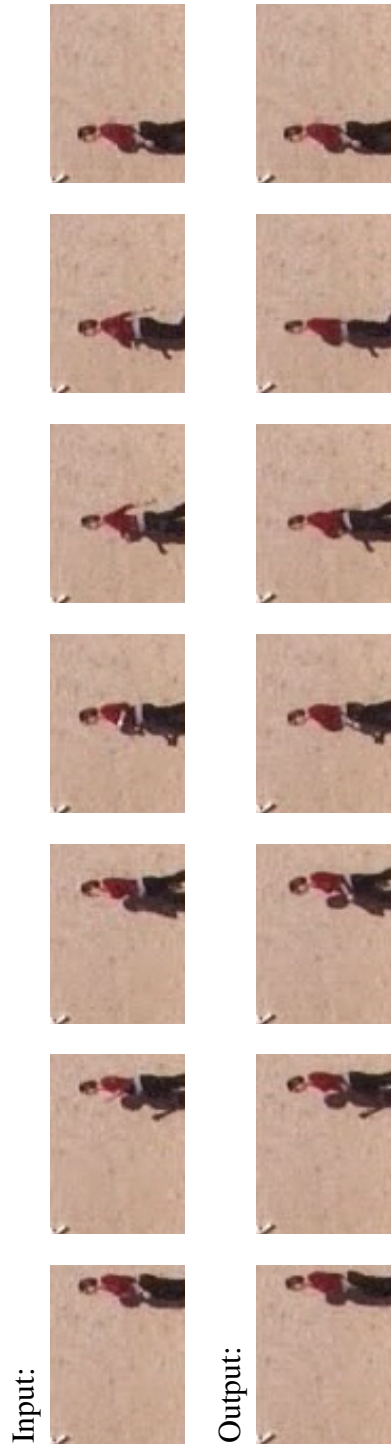


Fig. 15

MODIFICATION OF A VISUAL STORY. THE WOMAN IN THE TOP ROW SCRATCHES HER EAR BUT IN THE BOTTOM ROW SHE DOESN'T. THE VIDEO SIZE IS $90 \times 360 \times 190$ WITH 70,305 MISSING PIXELS.

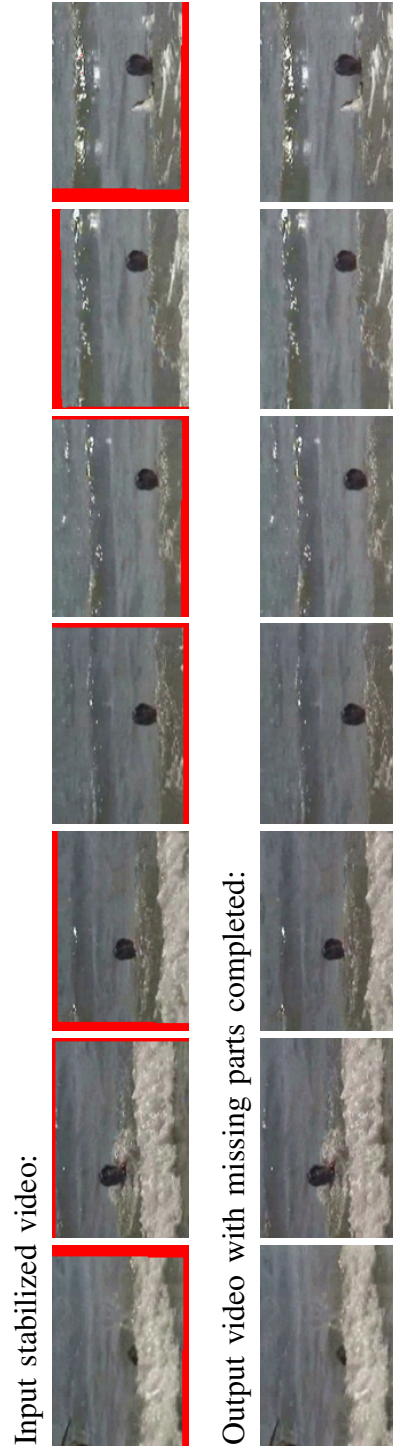


Fig. 16

STABILIZATION EXAMPLE. THE TOP ROW SHOWS A STABILIZED VIDEO SEQUENCE (COURTESY OF E. OFEK, [23]) OF SIZE $240 \times 352 \times 91$ WITH 466,501 MISSING PIXELS. EACH FRAME WAS TRANSFORMED TO CANCEL THE CAMERA MOTION AND AS A RESULT, PARTS OF EACH FRAME ARE MISSING (SHOWN HERE IN RED). TRADITIONALLY, THE VIDEO WOULD BE CROPPED TO AVOID THESE AREAS AND TO KEEP ONLY THE REGION THAT IS VISIBLE THROUGHOUT. HERE THE MISSING PARTS ARE FILLED UP USING OUR METHOD TO CREATE VIDEO WITH THE FULL SPATIAL EXTENT. THE RESULTS HERE ARE COMPARABLE WITH THOSE IN [23].

3. Modify a visual story: Our method can be used to make people in a movie change their behavior. For example, if an actor has absent-mindedly picked his nose during a film recording, then the video parts containing the obscene behavior can be removed, to be coherently replaced by information from data containing a range of “acceptable” behaviors. This is demonstrated in Fig. 15 where an unwanted scratching of the ear is removed.

4. Complete field-of-view of a stabilized video: When a video sequence is stabilized, There will be missing parts in the perimeter of each frame. Since the hole does not have to be bounded, the method can be applied here as well. Figure 16 shows such an application.

5. Creation of video textures: The method is also capable in creating large video textures from small ones. In Fig. 17, a small video sequence (32 frames) was extended to a larger one, both spatially and temporally.

VII. SUMMARY AND CONCLUSIONS

We have presented an objective function for completion of missing data and an algorithm to optimize it. The objective treats the data uniformly in all dimensions and the algorithm was demonstrated on the completion of several challenging video sequences. In this realm, the objective proved to be very suitable as it only relies on very small, local parts of information. It successfully solved problems with hundreds of thousands of unknowns. The algorithm takes advantage of the sparsity of the database within the huge space of all image patches to derive an update step.

The method provides a principled approach with a clear objective function that extends those used in other works. We have shown here that it can be seen as a variant of EM and we have also shown its relation to BP.

There are several possible improvements to this method. First, we did not attempt to prune the database in any way, even though this is clearly the bottleneck w.r.t. running time. Second, we applied the measure to windows around every pixel, rather than a sparser set. Using a sparser grid (say, every 3rd pixel) will give a significant speedup (3^3). Third, breaking the database into meaningful portions can prevent the algorithm to from completing one class of data with another, even if they are similar (e.g. complete a person with background), as was done in [18], [24]. Fourth, the proposed method does not attempt to coalesce identical pixels that come from

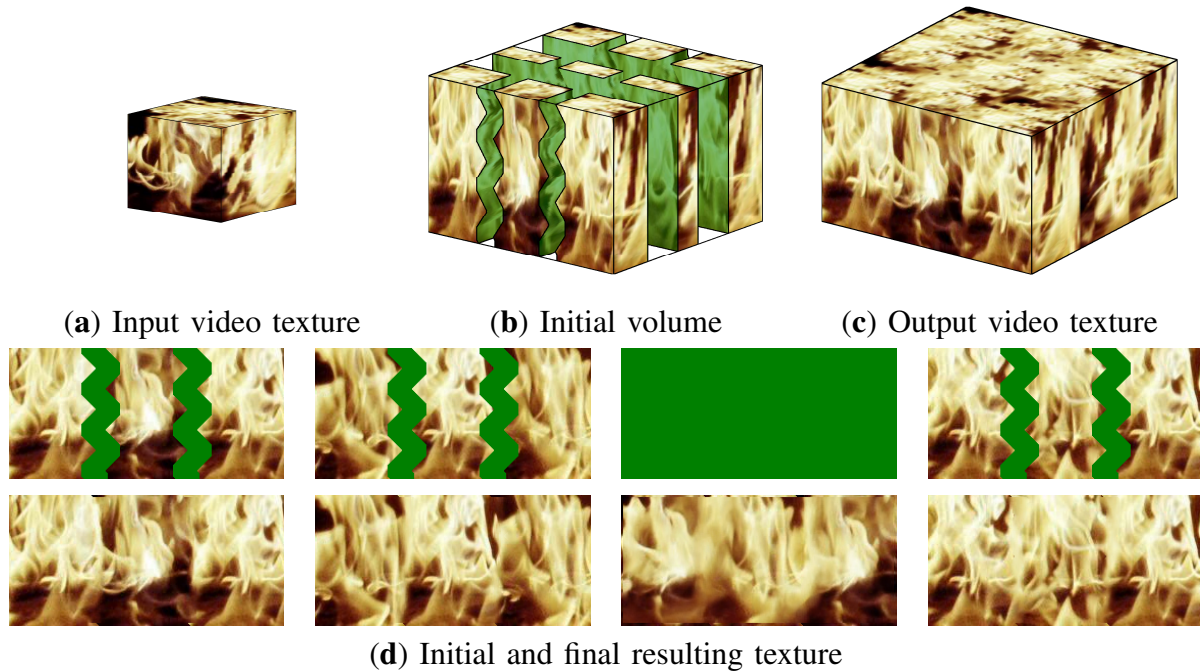


Fig. 17

VIDEO TEXTURE SYNTHESIS. (a) A SMALL VIDEO (SIZED $128 \times 128 \times 32$) IS EXPANDED TO A LARGER ONE, BOTH IN SPACE AND IN TIME. (b) IS THE LARGER VOLUME OF $128 \times 270 \times 60$ FRAMES WITH 1,081,584 MISSING PIXELS. SMALL PIECES OF (a) WERE RANDOMLY PLACED IN THE VOLUME WITH GAPS BETWEEN THEM. IN (c) THE GAPS WERE FILLED AUTOMATICALLY BY THE ALGORITHM. SEVERAL FRAMES FROM THE INITIAL AND RESULTING VIDEOS ARE SHOWN IN (d).

the same (static) part of the scene. These can be combined to form joint constraint rather than being solved almost independently in each frame as was done here.

VIII. ACKNOWLEDGMENTS

The authors would like to thank Oren Boiman for pointing out the graphical model that describes our algorithm, and for the anonymous reviewers for their helpful remarks.

REFERENCES

- [1] S. Arya and D. M. Mount. Approximate nearest neighbor queries in fixed dimensions. In *SODA '93: Proceedings of the fourth annual ACM-SIAM Symposium on Discrete algorithms*, pages 271–280, Philadelphia, PA, USA, 1993. Society for Industrial and Applied Mathematics.
- [2] M. Ashikhmin. Synthesizing natural textures. In *Proceedings of the 2001 symposium on Interactive 3D graphics*, pages 217–226, New York, NY, USA, 2001. ACM Press.
- [3] Z. Bar-Joseph, R. El-Yaniv, D. Lischinski, , and M. Werman. Texture mixing and texture movie synthesis using statistical learning. In *Transactions on Visualization and Computer Graphics*, volume 7, pages 120–135, Piscataway, NJ, USA, 2001. IEEE Educational Activities Department.
- [4] M. Bertalmío, A. L. Bertozzi, and G. Sapiro. Navier-stokes, fluid dynamics, and image and video inpainting. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 355–362, December 2001.
- [5] M. Bertalmio, G. Sapiro, V. Caselles, and C. Ballester. Image inpainting. In *SIGGRAPH '00: Computer Graphics Proceedings*, pages 417–424, New York, NY, USA, 2000. ACM Press/Addison-Wesley Publishing Co.
- [6] T. Chan, S. H. Kang, and J. Shen. Euler’s elastica and curvature based inpainting. *Journal on Applied Mathematics*, 63(2):564–592, 2002.
- [7] V. Cheung, B. J. Frey, and N. Jovic. Video epitomes. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 42–49, 2005.
- [8] D. Comaniciu and P. Meer. Mean shift: A robust approach toward feature space analysis. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 24(5):603–619, May 2002.
- [9] A. Criminisi, P. Pérez, and K. Toyama. Object removal by exemplar-based inpainting. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, pages 721–728, June 2003.
- [10] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society*, 39(B):1–38, 1977.
- [11] G. Doretto, A. Chiuso, Y. Wu, and S. Soatto. Dynamic textures. *IEEE International Conference on Computer Vision (ICCV)*, 51(2):91–109, 2003.
- [12] I. Drori, D. Cohen-Or, and H. Yeshurun. Fragment-based image completion. In *Transactions on Graphics*, volume 22, pages 303–312, New York, NY, USA, 2003. SIGGRAPH, ACM Press.
- [13] A. Efros and W.T. Freeman. Image quilting for texture synthesis and transfer. In *SIGGRAPH '01, Computer Graphics Proceedings*, pages 341–346, New York, NY, USA, 2001. ACM Press.
- [14] A. Efros and T. Leung. Texture synthesis by non-parametric sampling. *International Journal of Computer Vision (IJCV)*, 2:1033–1038, 1999.
- [15] A.W. Fitzgibbon, Y. Wexler, and A. Zisserman. Image-based rendering using image-based priors. *International Journal of Computer Vision (IJCV)*, 63(2):141–151, 2005.
- [16] W. T. Freeman, E. C. Pasztor, and O. T. Carmichael. Learning low-level vision. *International Journal of Computer Vision (IJCV)*, 40(1):25–47, 2000.
- [17] M. Irani and S. Peleg. Motion analysis for image enhancement: Resolution, occlusion, and transparency. *Journal of Visual Communication and Image Representation (JVCIR)*, 4:324–335, 1993.
- [18] J. Jia, T. P. Wu, Y. W. Tai, and C. K. Tang. Video repairing: Inference of foreground and background under severe occlusion. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 364–371, 2004.

- [19] A. Kokaram. Practical, unified, motion and missing data treatment in degraded video. *Journal of Mathematical Imaging and Vision*, 20(1-2):163–177, 2004.
- [20] A.C. Kokaram, B. Collis, and S. Robinson. Automated rig removal with bayesian motion interpolation. *IEEE Proceedings - Vision, Image and Signal Processing*, 152(4):407–414, August 2005.
- [21] V. Kwatra, A. Schödl, I. Essa, G. Turk, and A. Bobick. Graphcut textures: Image and video synthesis using graph cuts. In *Transactions on Graphics (TOG)*, volume 22, pages 277–286, New York, NY, USA, 2003. SIGGRAPH, ACM Press.
- [22] A. Levin, A. Zomet, and Y. Weiss. Learning how to inpaint from global image statistics. In *IEEE International Conference on Computer Vision (ICCV)*, pages 305–312, Washington DC, 2003. IEEE Computer Society.
- [23] Y. Matsushita, E. Ofek, X. Tang, and H.Y. Shum. Full-frame video stabilization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 50–57, 2005.
- [24] K.A. Patwardhan, G. Sapiro, and M. Bertalmio. Video inpainting of occluding and occluded objects. In *IEEE International Conference on Image Processing (ICIP)*, pages 69–72, September 2005.
- [25] A. Schödl and I. Essa. Controlled animation of video sprites. In *SCA '02: Proceedings of the First ACM Symposium on Computer Animation (held in Conjunction with ACM SIGGRAPH 2002)*, pages 121–127, New York, NY, USA, 2002. ACM Press.
- [26] A. Schödl, R. Szeliski, D. H. Salesin, and I. Essa. Video textures. In *SIGGRAPH '00, Computer Graphics Proceedings*, pages 489–498, New York, NY, USA, 2000. ACM Press/Addison-Wesley Publishing Co.
- [27] U. Trottenber, C. Oosterlee, and A. Schüller. *Multigrid*. Academic Press, Inc., Orlando, FL, USA, 2000.
- [28] L. Wei and M. Levoy. Fast texture synthesis using tree-structured vector quantization. In *SIGGRAPH '00: Computer Graphics Proceedings*, pages 479–488, New York, NY, USA, 2000. ACM Press/Addison-Wesley Publishing Co.
- [29] Y. Wexler, E. Shechtman, and M. Irani. Space-time video completion. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 120–127, 2004.

APPENDIX

In this section we derive the optimization algorithm from section III using the statistical estimation framework for the interested reader. In this framework, the global visual coherence of Eq. (1) can be derived as a likelihood function via a graphical model, and our iterative optimization process can be viewed as a variant of the EM algorithm to find the maximum likelihood estimate.

Under this model the unknown parameters, denoted by Θ , are the color values of the missing pixel locations p in the space-time hole: $\Theta = \{c_p | p \in \mathcal{H}\}$. The known pixels around the hole are the observed variables $\mathbf{Y} = \{c_p | p \in \mathcal{S} \setminus \mathcal{H}\}$. The windows $\{W_n\}_1^N$ are defined as small space-time patches overlapping the hole where N is their total number. The set of patches $\mathbf{X} = \{X_n\}_1^N$ are the hidden variables corresponding to W_n . The space of possible assignments for each X_n are the dataset patches in \mathcal{D} . We assume that each X_n can have any assignment with some probability. For the completion examples here, the dataset is composed of all the patches from the same sequence that are completely known, i.e. $\mathbf{X} = \{X_n | X_n \subset \mathcal{S} \setminus \mathcal{H}\}$ but the dataset may also contain patches from another source.

This setup can be described as a graphical model which is illustrated in Fig. 18 in one dimension. Each patch W_n is associated with one hidden variable X_n . This in turn is connected

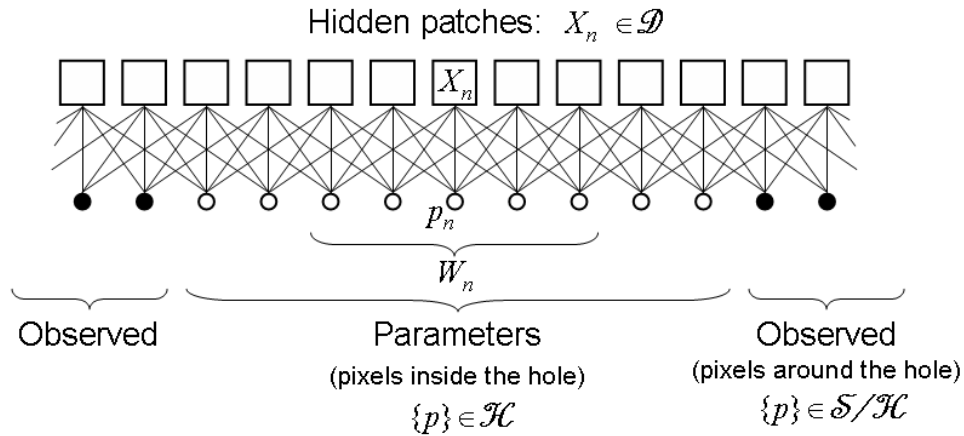


Fig. 18

OUR GRAPHICAL MODEL FOR THE COMPLETION PROBLEM

to all the pixel locations it covers, $p \in W_n$. As each W_n denotes a patch overlapping the hole, at least some of these pixel locations are unknowns, i.e. belong to Θ . Let $c_p = W_n^p$ denote the color of the pixel p in the appropriate location within the window W_n and let X_n^p denote the color at the same location within the dataset windows. Note that while the color values X_n^p corresponding to pixel p may vary for different window locations n , the actual pixel color c_p is the *same* for all overlapping windows W_n^p . Using these notations, an edge in the graph that connects an unknown pixel p with an overlapping window X_n has an edge potential of the form $\phi(c_p, X_n) = e^{-\frac{1}{2\sigma^2}(c_p - X_n^p)^2}$.

The graph in Fig18 with the definition of its edge potentials is equivalent to the following joint probability density of the observed boundary variables and the hidden patches given the parameters Θ :

$$f(\mathbf{Y}, \mathbf{X}; \Theta) = \beta \prod_{n=1}^N \prod_{p \in W_n} \phi(c_p, X_n) = \beta \prod_{n=1}^N \prod_{p \in W_n} e^{-\frac{(c_p - X_n^p)^2}{2\sigma^2}}$$

where p denotes here either missing or observed pixel and β is a constant normalizing the product to 1. This product is exactly the similarity measure defined previously in Eq.(2):

$$f(\mathbf{Y}, \mathbf{X}; \Theta) = \beta \prod_{n=1}^N e^{-d(X_n, W_n)} = \beta \prod_{n=1}^N \text{sim}(X_n, W_n)$$

To obtain the likelihood function we need to marginalize over the hidden variables. Thus we need to integrate over all possible assignments for the hidden variables $X_1 \dots X_N$:

$$L = f(\mathbf{Y}; \Theta) = \sum_{(X_1, \dots, X_N) \in \mathcal{D}^N} f(\mathbf{Y}, \mathbf{X}; \Theta) \quad (5)$$

$$= \sum_{(X_1, \dots, X_N) \in \mathcal{D}^N} \beta \prod_{n=1}^N \text{sim}(X_n, W_n) \quad (6)$$

The maximum likelihood solution to the completion problem under the above model assumptions is the set of hole pixel values (Θ) that maximize this likelihood function. Note that since the summation terms are products of all patch similarities, we will assume that this sum is dominated by the maximal product value (deviation of one of the patches from its best match will cause a significant decrease of the entire product term): $\max L \approx \max_{\{\mathbf{X}\}} \beta \prod_{n=1}^N \text{sim}(X_n, W_n)$. Given the point values, seeking the best patch matches can be done independently for each W_n , hence we can change the order of the max and product operators:

$$\max L \approx \beta \prod_{n=1}^N \max_{X_n} \text{sim}(X_n, W_n)$$

meaning that the *maximum likelihood* solution is the same completion that attains best *visual coherence* according to Eq. 1.

We will next show that our optimization algorithm fits the EM algorithm [10] for maximizing the above likelihood function.

In the **E** step, at iteration t the posterior probability density $\hat{\pi}^t$ is computed. Due to conditional independence of the hidden patches, this posterior can be written as the following product:

$$\hat{\pi}^t = f(\mathbf{X}|\mathbf{Y}; \Theta^t) = \prod_{n=1}^N f(X_n|\mathbf{Y}; \Theta^t) = \prod_{n=1}^N \beta_n \text{sim}(X_n, W_n)$$

Thus we get a probability value $\hat{\pi}_t$ for each possible assignment of the dataset patches. Given the above considerations, these probabilities vanish in all but the best match assignments in each pixel thus the resulting PDF is an indicator function. This common assumption, also known as ‘‘Hard EM’’, justifies the choice of one nearest neighbor.

In the **M** step, the current set of parameters Θ are estimated by maximizing the following function:

$$\hat{\Theta}^{t+1} = \arg \max_{\Theta} \left(\sum_{(X_1, \dots, X_N) \in \mathcal{D}^N} \hat{\pi}^t \log f(\mathbf{X}, \mathbf{Y}; \Theta) \right)$$

Given the graphical model in Fig.18, each unknown p depends only on its overlapping patches and the pixels are conditionally independent. Thus the global maximization may be separated to local operations on each pixel:

$$\hat{c}_p^{t+1} = \arg \max_{c_p} \left(- \sum_{\{n \mid p \in W_n\}} (c_p - \hat{X}_n^p)^2 \right)$$

where the \hat{X}_n are the best patch assignments for the current iteration (denoted also by V in Sec. III-A). This is an L_2 distance between the point value and the corresponding color values in all covering patches and it is maximized by the *mean* of these values similar to Eq.(3).

Similar to the classic EM, the likelihood in the ‘‘Hard EM’’ presented here is increased in each **E** and **M** steps. Thus the algorithm converges to some local maxima. The use of a multi-scale process (see Sec. III-B) and other adjustments (see Sec. III-A) leads to a quick convergence, and to a realistic solution with high likelihood.