

Spectral Total-Variation Local Scale Signatures for Image Manipulation and Fusion

Ester Hait and Guy Gilboa, *Senior Member, IEEE*

Abstract—We propose a unified framework for isolating, comparing and differentiating objects within an image. We rely on the recently proposed total-variation transform, yielding a continuous, multi-scale, fully edge-preserving, local descriptor, referred to as spectral total-variation local scale signatures. We show and analyze several useful merits of this framework. Signatures are sensitive to size, local contrast and composition of structures; are invariant to translation, rotation, flip and linear illumination changes; and texture signatures are robust to the underlying structures. We prove exact conditions in the 1D case.

We propose several applications for this framework: saliency map extraction for fusion of thermal and optical images or for medical imaging, clustering of vein-like features and size-based image manipulation.

Index Terms—Spectral Total-Variation, Image Fusion, Image Segmentation, Edge Detection, Size Differentiation, Clustering, Saliency, Thermal Imagery, Medical Imagery.

I. INTRODUCTION

DIFFERENTIATING objects within an image by contrast, size or structure is a fundamental image processing task. It is highly useful for various image modalities and applications, such as image clustering, enhancement and fusion. A key feature in many modalities (natural images, medical, thermal, depth etc.) are edges, or discontinuities in the data.

For this purpose, the spectral total variation transform (spectral TV) has been recently introduced as a useful edge-preserving, multi-scale decomposition tool [1],[2],[3],[4]. Some previous spectral TV-based approaches succeeded in texture extraction and manipulation [1],[2] or texture-structure decomposition [3],[5], while other TV or spectral TV-based approaches successfully differentiated objects by size [4],[6],[7]. However, no previous approach is suited for the task of object differentiation by both size and contrast (Fig. 3). Applying the transform to different modalities is also challenging, due to their complex nature, multi-scaled and occluded content.

In this paper, we present a novel, unified framework for object differentiation by contrast, size or structure, including complex multi-scaled objects. We capture salient objects by exploiting the comprehensive scale and location information extracted from spectral TV, referred to as *Spectral TV Local Scale Signatures* (Fig. 1). Stemming from an edge preserving, sparse spectral transform, signatures of significant objects are sparse and strong; their locality allows good differentiation within an image. We show and analyze the essential properties of the signatures: sensitivity to size, local contrast and

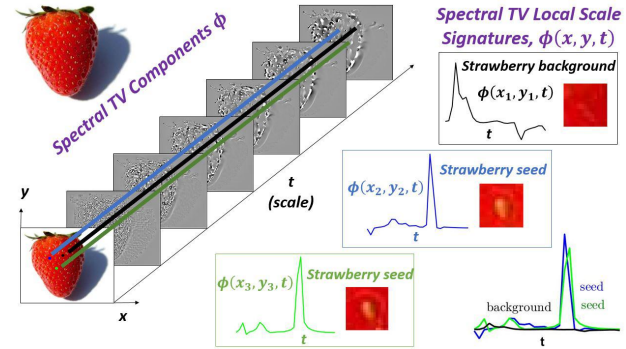


Figure 1: **Spectral TV Local Scale Signatures**, $\phi(x, y, t)$: a multi-scale spectral TV per-pixel description, sensitive to size, local contrast and composition of structures, invariant to translation, rotation, flip and linear illumination change, with texture invariance to underlying structure. Pixels with common features (strawberry seeds) have similar signatures; different pixels can be differentiated by their distinct signatures.

composition of structures; invariance to rotation, translation, flip and linear illumination change; and invariance of texture to structure. We show how no previous spectral TV-based approach can handle object differentiation by size and contrast, as well as a composition of structures (Fig. 3).

Though relying merely on a few simple cornerstones, our algorithm is applicable to various tasks (Fig. 2): fusion of thermal and visible images, or of medical images of different modalities; clustering of vein-like repetitive structures (segmentation / edge detection); and size differentiation.

The rest of the paper is organized as follows. Sec. II briefly surveys previous image descriptors. Sec. III introduces spectral TV (Sec. III-A), and reviews and analyzes previous approaches. Sec. IV presents the novel concept of spectral TV signatures and their properties. Sec. V gives a theoretical analysis. Sec. VI presents a unified framework for object differentiation, demonstrated for synthetic images, and applied for image manipulation (VI-C) and fusion (VI-D). Sec. VII gives experimental results and comparisons to other methods. Sec. VIII concludes our work. The Appendix provides further analysis and suggests new fusion visualization methods.

II. RELATED WORK

In the past decades, high-level image understanding and processing has considerably relied on feature descriptors of different types. For example, sparse descriptors, describing characteristics of interest points only (e.g. corners) [8], [9].

E. Hait and G. Gilboa are with the Department of Electrical Engineering, Technion - Israel Institute of Technology, Haifa, Israel, e-mail: ety-hait@campus.technion.ac.il, guy.gilboa@ee.technion.ac.il

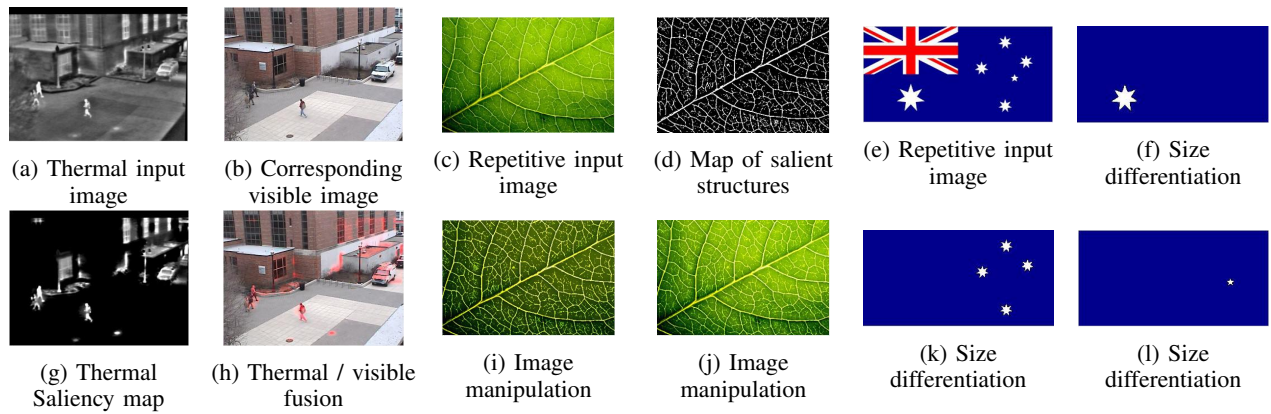


Figure 2: Image fusion (a-b,g-h) and image manipulation (c-f,i-l) using spectral TV local scale signatures.

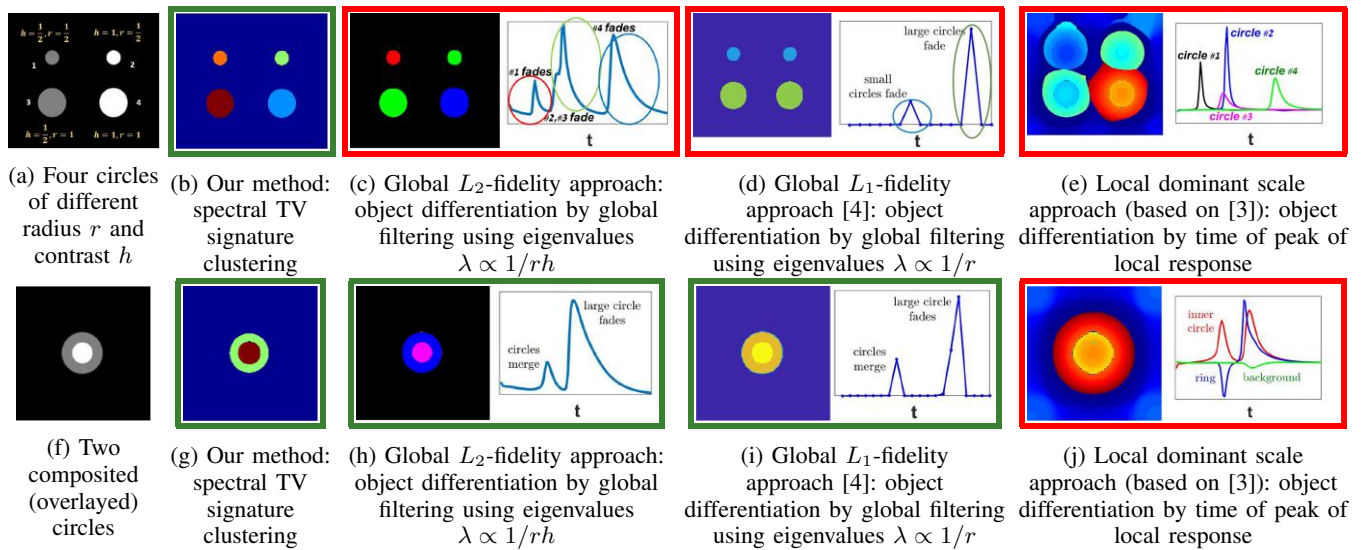


Figure 3: **Top row: no spectral-TV based method but ours (b) can differentiate objects by both size and contrast.** Global filtering (c,d): different objects have the same eigenvalue (inversely proportional to scale); they thus respond simultaneously in a spectrum single peak, and cannot be differentiated. Local dominant scale (e): pixels of different objects have similar dominant scales, and thus cannot be differentiated. Bottom row: our method can also differentiate composited (overlaid) structures.

Our focus here, however, is on dense descriptors, describing every image pixel using properties of the pixel and its surroundings. Various algorithms [6], [10], [11], [12] rely on fundamental properties, such as intensity and color, gradient magnitude and orientation, textures and patterns. Some of these algorithms use features for learning [13], [14]. Among these descriptors, those based on histograms of patterns [15] or of oriented gradients [16] are rotation-variant, thus suitable for texture applications. Other approaches include multi-scale decompositions using image transforms [17] or diffusion [18]; or using the responses to a set of predefined linear filters as features [19], [20], [21]. Recently, convolutional neural networks (CNN) have emerged as a highly successful data-driven tool [22], [23], [24]. However, this feature learning approach is not generic and lacks solid theoretical background. To overcome this, recent CNN methods use sets of linear filters in their first convolutional layers [25], [26].

In this paper, we present a novel dense, multi-scale, edge preserving, local descriptor. Its properties are highly suited for

detecting and clustering vein- and disk-like structures and for constructing saliency maps for fusion.

III. PRELIMINARIES

A. Spectral Total Variation

The total variation (TV) functional [27] has been widely used as an image regularizer, e.g. for denoising and deconvolution [27], [28], [29], [30], decomposition and texture analysis [31], [32], [33] and fusion [34], [35]. Denoting the image domain as Ω , and the gradient (understood as the distributional gradient) as ∇ , the TV functional is defined as:

$$J(u) = \int_{\Omega} |\nabla u(x)| dx. \quad (1)$$

For an input image $f(x)$, the gradient descent evolution of this functional, known as TV flow [36], is defined as:

$$\begin{aligned} \frac{\partial u}{\partial t} &= \operatorname{div} \left(\frac{\nabla u}{|\nabla u|} \right) \quad \text{in } (0, \infty) \times \Omega, \\ u(x, 0) &= f(x) \quad \text{in } x \in \Omega, \end{aligned} \quad (2)$$

with Neumann boundary conditions (note that the right-hand-side of (2) should be understood as a negative subgradient of $J(u)$). Performing TV flow, up to a scale T , is analogous to solving the ROF [27] optimization problem [1]:

$$\arg \min_u \left\{ J(u) + \frac{1}{2T} \|f - u\|_2^2 \right\}. \quad (3)$$

Denoting $t \in [0, \infty)$ as the **time** or **scale** parameter of the TV flow (2), and u_{tt} as the second time derivative of its solution $u(x, t)$, the spectral TV transform [37] is defined by:

$$\phi(x, t) = u_{tt}(x, t)t. \quad (4)$$

We refer to $\phi(x, t)$ as the spectral component or band at scale t (see Fig. 1, left for a visualization of ϕ at some scales). The main properties of this nonlinear transform are as follows:

- 1) **Reconstruction.** Any zero mean, bounded variation function f can be reconstructed by $f(x) = \int_0^\infty \phi(x, t) dt$.
- 2) **Spectral representation.** We define a nonlinear eigenfunction with respect to the subdifferential of TV, $\partial J(u)$, as: $\lambda u \in \partial J(u)$. Examples of such eigenfunctions are disks, or convex sets with low curvature on the boundary, bounded by the perimeter to area ratio [36]. Then the spectral response of an eigenfunction f with **eigenvalue** λ is a Dirac delta at **scale** $1/\lambda$: $\phi(t) = \delta(t - 1/\lambda)f$. Thus, a signal composed of spatially-separated eigenfunctions has a highly sparse representation (analogous to Fourier transform representation of sine functions).
- 3) **Orthogonality.** Under some conditions, such as in the discrete 1D case, spectral components are orthogonal to each other: $\langle \phi(t_1), \phi(t_2) \rangle = 0, \forall t_1 \neq t_2$ [1].
- 4) **Filtering.** Given a filter $H(t)$, extending the reconstruction formula, $f_H(x) = \int_0^\infty H(t)\phi(x, t)dt$, allows the design of various edge-preserving TV filters.
- 5) **Translation and rotation invariance.** The ϕ components inherit the properties of the TV functional and are translation and rotation invariant.

These properties apply to the multiscale representation $\phi(x_0, t)$ of each pixel x_0 . This pixel-neighborhood relation characterization can successfully serve as a generic pixel descriptor.

We later use two additional definitions. First, the spectrum $S^f(t)$ is defined (using the original definition of [37]) as:

$$S^f(t) = \|\phi(x, t)\|_{L^1} = \int_{\Omega} |\phi(x, t)| dx, \quad (5)$$

which can be seen as the L^1 amplitude of the response at each scale $t \in [0, \infty)$. Second, the residual image $f_r(x)$, generated after some *finite* time T , is defined as:

$$f_r(x, T) = u(x, T) - u_t(x, T) \cdot T. \quad (6)$$

Further discussions can be found in [37].

B. Previous Spectral TV-based Approaches

Most previous spectral TV-based methods perform *global* scale analysis, that is, spatially integrating spectral information for each scale. Global spectral methods using the TV-flow (or L_2 -fidelity as in (3)) have been successfully used for texture

extraction and manipulation [1], [2], [5]. However, differentiating objects *within* an image is challenging, as different objects may correspond similarly in the global spectrum (5) (although the full signature is different).

Another successful global approach uses an L_1 -fidelity term in (3) for *contrast-invariant* multi-scale decomposition, differentiating objects only by their size [4]. However, we are interested in differentiating objects both by size and contrast. The only previous *local* approach uses local dominant scale analysis for structure-texture decomposition [3]. It relies on the time of the response peak of each pixel to fit a separation surface between two spectral bands: texture and structure. Despite its success for texture-structure decomposition, it is not suited for object differentiation. Despite its simplicity, it is limited, as it uses only the response *peak* information.

Toy Example Analysis: Fig. 3 studies the limitations of previous spectral TV-based methods.

Objects of various sizes and contrasts (Fig. 3a) cannot all be differentiated. In the global L_2 -fidelity method, two objects, one double the size, but half the contrast of the other, have the same eigenvalue $\lambda \propto 1/rh$, inversely proportional to scale. Thus, both respond simultaneously, forming a mutual peak (Fig. 3c) in the global spectrum (5). The global L_1 -fidelity method [4] cannot differentiate same-size, different-contrast objects, which have the same eigenvalue $\lambda \propto 1/r$. Thus, these objects present same-scale, simultaneous spectrum responses (Fig. 3d). Differentiating using the *local* (pixel-wise) dominant scale (following [3]) also fails, as different objects and background regions respond simultaneously (Fig. 3e).

A composition of structures (Fig. 3f) simulates real multi-scaled images. Composited objects impact each other's behavior, first merging, then fading. Both global methods can approximately differentiate objects (Figs. 3h, 3i). However, differentiating by local dominant scale fails (Fig. 3j).

In conclusion, no previous spectral-TV based method can differentiate objects by size and contrast, as well as a composition of structures. Only our method succeeds in this task (Figs. 3b, 3g), by exploiting more spectral TV information within a local framework.

IV. SPECTRAL TV LOCAL SCALE SIGNATURES

Object differentiation requires exploiting detailed, local, multi-scale information to handle different sizes, contrasts and complex structures. We thus introduce the concept of spectral TV local scale signatures. We denote the signatures of a signal $f(x)$ as $\phi_f(x, t)$, where ϕ is defined by (4). For each pixel there exists a well-defined, unique representation in the scale continuum (unlike classical pyramidal multiscale representations and wavelets), yielding a natural multi-scale pixel descriptor. This section introduces the properties of signatures and their implications for general images, including visual demonstrations. The interested reader can find an elaborated theoretical analysis in Sec. V.

A. Properties of Spectral TV Signatures

We now summarize the main properties of spectral TV signatures and illustrate them graphically by simple toy examples.

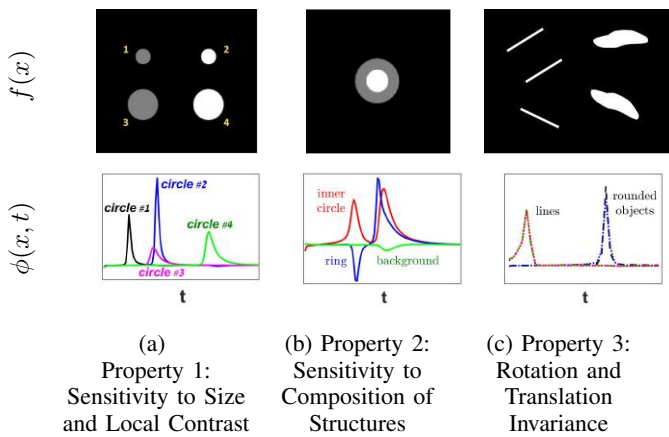


Figure 4: 2D Demonstration of properties 1, 2, 3. Signatures are distinct due to their sensitivity to size and contrast (a), and their sensitivity to composition of structures (b). However, they are invariant to rotation and translation (c).

Sensitivity Properties of Signatures:

Property 1. Sensitivity to Size and to Local Contrast

Signatures are sensitive to change in size (spatial scaling by factor a) and change in local contrast (contrast change by factor a) by the following relations:

$$\begin{aligned}\phi_{f(ax)} &= a\phi_f(ax, at), \\ \phi_{af(x)} &= \phi_f(x, t/a).\end{aligned}\quad (7)$$

This results from spectral TV properties [37]. This allows differentiating objects by their distinct signatures (Fig. 4a).

Property 2. Sensitivity to Composition of Structures

Structures of comparable scales respond differently when composed together (with spatial overlay), compared to their individual responses. That is, in this case the non-linearity of the TV transform applies, allowing for some signals f, g :

$$\phi_{f+g} \neq \phi_f + \phi_g. \quad (8)$$

Composited objects thus have distinct signatures (Fig. 4b). An analytical solution (Appendix A, Proposition 1) and a demonstration (Fig. 5) are given for the 1D staircase signal. Note, that non-linearity is marginal for structures of very different scales or for spatially distant ones.

Invariance Properties of Signatures:

Property 3. Invariance to Rotation, Translation and Flip

Signatures are rotation and translation invariant in \mathbb{R}^n :

$$\begin{aligned}\phi_{f(Rx)} &= \phi_f(Rx, t), \\ \phi_{f(x-a)} &= \phi_f(x-a, t),\end{aligned}\quad (9)$$

where $R(x)$ is a rotation matrix, and a is a spatial shift (Fig. 4c). This also results from spectral TV properties [37]. Since TV is invariant to the coordinate system (being rotationally invariant and sensitive only to derivatives), signatures are also invariant to flip w.r.t. an arbitrary coordinate system:

$$\phi_{f(x)}(x) = \phi_{f(-x)}(-x). \quad (10)$$

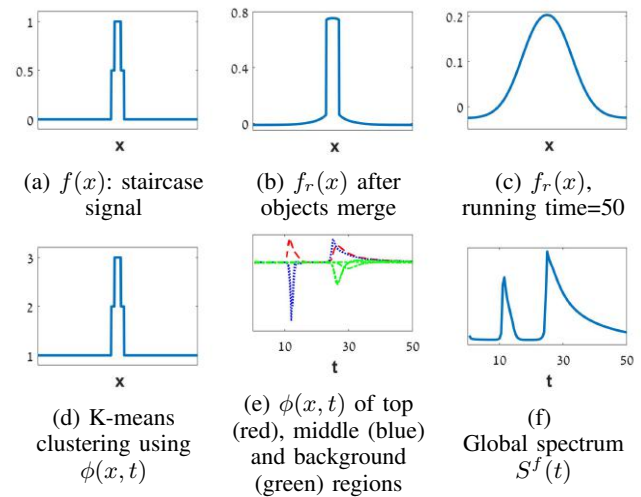


Figure 5: Demonstration of Property 2 for the 1D staircase signal (a). Spectral TV signatures, generated during TV flow of the signal (b,c), are distinct for different regions (e). Thus, signature clustering allows differentiating the regions (d).

Fig. 7 shows how these properties can be useful for finding similar image textures, where a patch-based comparison fails.

Property 4. Invariance of Texture to Structure

Signatures of textures (patterns) are invariant to their underlying structures for fine scales under very broad conditions (precise conditions are given in Sec. V-C). Fig. 6 (top) shows the invariance of signatures of objects with identical textures to their different underlying global contrasts for fine scales. A 1D proof (Sec. V-C, Theorem 1) and a 1D demonstration (Fig. 9, top) are given.

Property 5. Invariance to Linear Illumination Change

Signatures are invariant to a linear change of illumination, which holds up to a certain scale (precise conditions are given in Sec. V-D). Fig. 6 (bottom) shows the invariance of signatures of objects with identical textures and identical global contrasts. A 1D proof (Sec. V-D, Theorem 2) and a 1D demonstration (Fig. 9, bottom) are given.

B. Implications for General Images

Fig. 8 shows examples of signature similarity and distinctness for different modalities. Comparing the signatures of different-size, same-color objects in Fig. 8a, signatures of small white stars display stronger, earlier responses, distinct from those of the larger white stripes (Fig. 8d). For thermal and medical images (Fig. 8b), signatures of objects different in size or contrast are distinct (Fig. 8c). Moreover, signatures of highly-contrasted objects display stronger responses than those of weakly-contrasted ones (Fig. 8e). Signature enhancement (see Sec. VI-D) improves the distinctness of signatures of different groups (Fig. 8f). This is useful since the salient objects in these modalities are usually highly-contrasted.

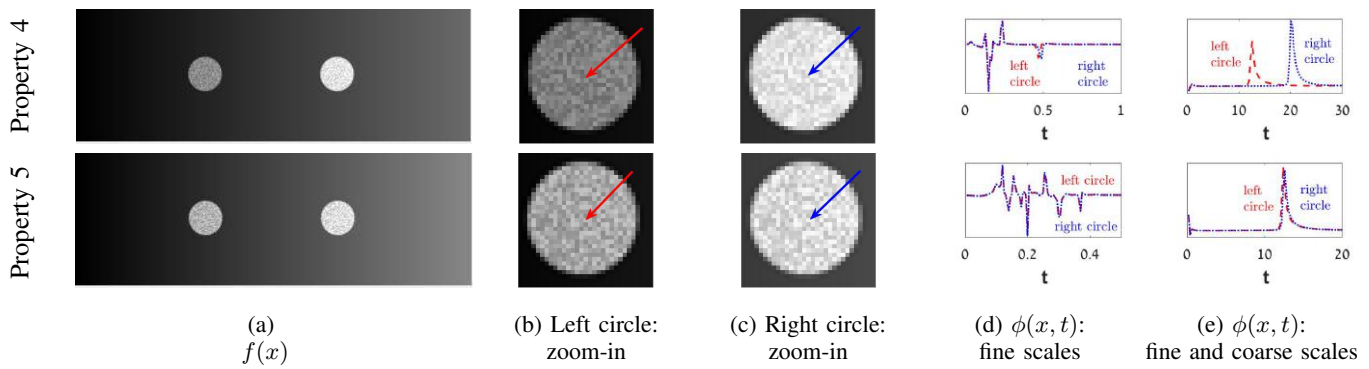


Figure 6: 2D Demonstration of properties 4,5. Signatures of objects with **identical textures (patterns)** are invariant to their different underlying structures - in this case, their different global contrasts - for fine scales (top). Signatures of objects with **identical textures and identical structures** are identical - thanks to their invariance to linear illumination changes (bottom).

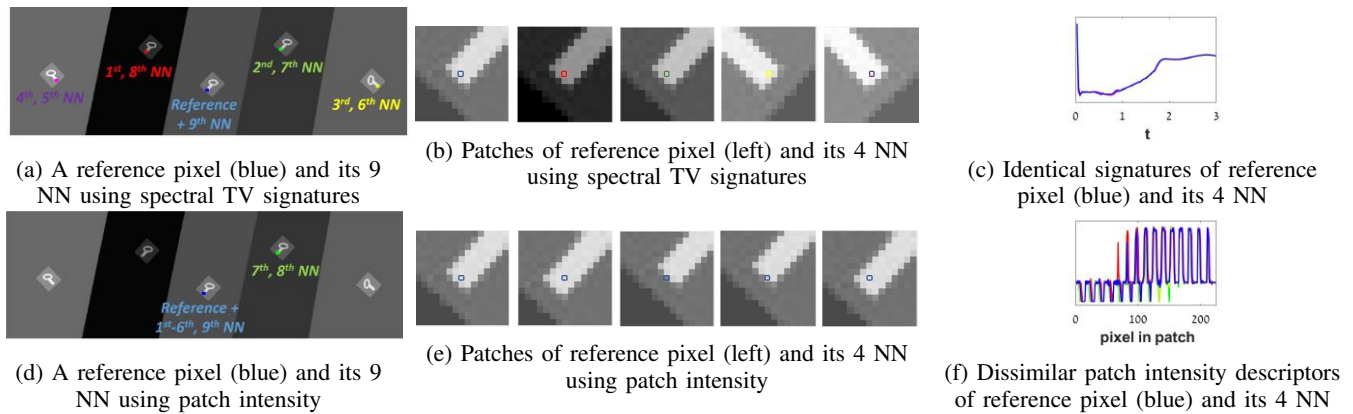


Figure 7: 2D Demonstration of the invariance to rotation, translation and flip, and the invariance of textures to their underlying structures. Using spectral TV signatures, the 9 nearest neighbors (NN) of a reference pixel (reference marked in blue) are pixels of similar textures, but of different global contrasts, rotations, translations or flips (top). Conversely, the patch intensity descriptor fails to find these texturally-similar pixels (bottom).

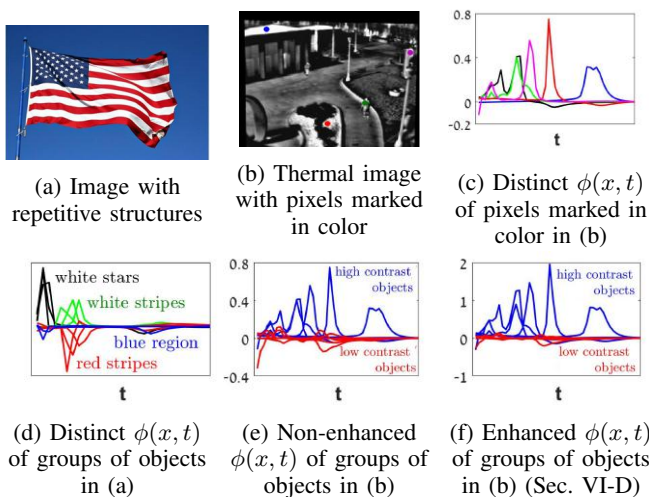


Figure 8: Distinctness of signatures of groups of objects with common features. For an image with repetitive structures (a), groups of objects have distinct signatures (d). For a thermal image, pixels of different objects (marked in color, b) have distinct signatures (c), and so do groups of objects (e,f).

V. THEORETICAL ANALYSIS

This section presents theoretical analysis of some signature properties in the 1D case for the theory-oriented reader. We formulate a sufficient condition for local patterns to merge first; and give sufficient conditions and proofs of the invariance to linear illumination change and of texture to structure. Appendix A shows an analytic solution of the separation of composited regions for the 1D staircase signal.

A. Preliminaries

Our analysis below is based on the work of Steidl *et al.* [38], which gives an analytic solution for the TV-flow in the time continuous, spatially discrete 1D case. They show that in a TV-flow evolution, each pixel belongs to a local constant region (all pixels in the region are connected and have the same value), which dictates its behavior. The region evolves at a certain constant speed (inversely proportional to region size), until a *merging event* occurs, that is, two neighboring regions obtain the same value. Let $f \in \mathbb{R}^N$ be a discrete 1D input signal of size N pixels. Let $u \in \mathbb{R}^N \times [0, \infty)$ be the space-discrete realization of the TV-flow, defined by (2). We

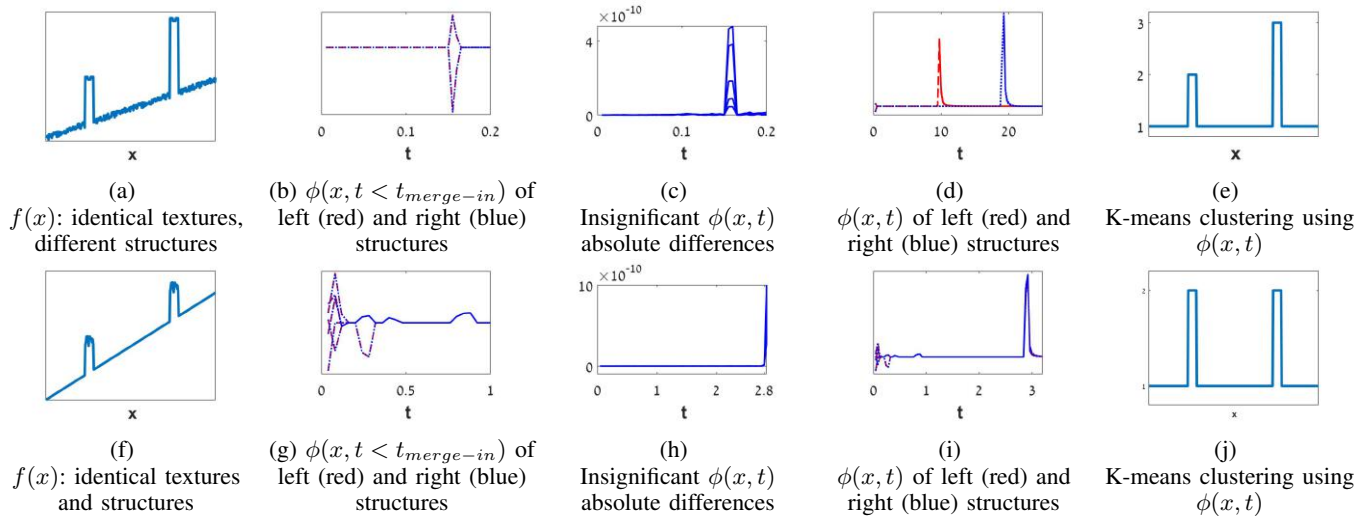


Figure 9: Demonstration of Property 4 (top) and Property 5 (bottom) for 1D signals. A signal with identical textures, but different underlying structures (a) has identical fine-scale signatures - up to $t_{merge-in}$, when local patterns merge (b,c). Signatures then become distinct (d). A signal with underlying structures of identical global contrasts (f), not only has identical signatures for $t < t_{merge-in}$ (g), but also identical signatures up to $t_{merge-struct}$, when structures merge with background simultaneously (h,i). Thus, signature clustering allows differentiating between structures (e), or structures from background (j).

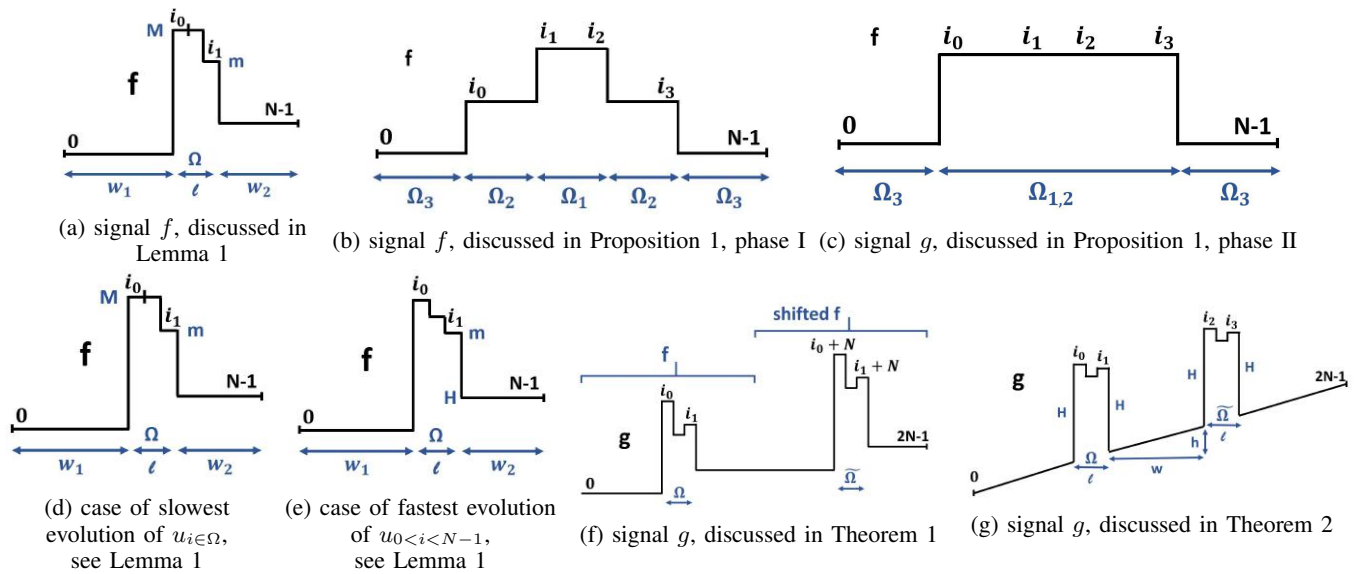


Figure 10: Signals discussed in Lemma 1 (a,d,e), Theorem 1 (f), Theorem 2 (g) and Proposition 1 (b,c).

denote by $u_i(t)$ the value of u at pixel i at time t . Two main properties of this dynamic are:

- 1) There exists a finite number of merging events, $0 = t_0 < t_1 < \dots < t_{n-1} < t_n$ (Proposition 4.1 (ii) in [38]).
- 2) Within the time intervals between merging events, $t \in [t_j, t_{j+1})$, all pixels u_i , belonging to a constant region $\{u_{i-l+1}, \dots, u_{i+r}\}$ of size w_{i,t_j} with relation μ_{i,t_j} to its neighboring regions, evolve linearly (4.1 (iii) in [38]):

$$u_i(t) = u_i(t_j) + \mu_{i,t_j} \frac{2(t-t_j)}{w_{i,t_j}},$$

$$\mu_{i,t_j} = \begin{cases} 0, & \text{if } \{u_{i-l}, \dots, u_{i+r+1}\} \text{ is strictly monotonic} \\ 1, & \text{if } u_i \text{ is minimal in } \{u_{i-l}, \dots, u_{i+r+1}\} \\ -1, & \text{if } u_i \text{ is maximal in } \{u_{i-l}, \dots, u_{i+r+1}\}. \end{cases} \quad (11)$$

B. Local Patterns Merge First

This section gives a 1D proof that regions of local patterns merge first, and only then merge with their surroundings. Let $f : \{0, \dots, N-1\} \rightarrow \mathbb{R}$ be as depicted in Fig. 10a, and let $\Omega_1 = \{0, \dots, i_0-1\}$, $\Omega = \{i_0, \dots, i_1\}$, $\Omega_2 = \{i_1+1, \dots, N-1\}$, where Ω_1, Ω_2 are constant regions (no patterns outside Ω). We

define: $w_1 \triangleq i_0$, $l \triangleq i_1 - i_0 + 1$, $w_2 \triangleq N - 1 - i_1$, and assume, without loss of generality, that $l < w_2 < w_1$ and $f[i_1 + 1] > f[i_0 - 1]$. We define: $m \triangleq \min_{i \in \Omega} f_i$ is attained at i_{min} , $M \triangleq \max_{i \in \Omega} f_i$ is attained at i_{max} , $H \triangleq f[i_1 + 1]$. We also define the following two **critical time points**: $t_{merge-in}$, the maximal merging time of Ω , and $t_{merge-out}$, the minimal merging time of $\{i_0, \dots, N - 1\}$.

Lemma 1 (Local Patterns Merge First). *Let f be as defined above. If $\frac{M-m}{m-H} \leq \frac{w_2}{l}$, then $t_{merge-in} < t_{merge-out}$.*

Proof. Relying on Section V-A, the TV flow dynamics of pixel i are determined by w_i , μ_i , regardless of $u_i(t = 0)$. The key concept of our proof is that μ_i depends only on pixels at the immediate edge of the region. **Thus, pixel behavior is not influenced by pixels "beyond" the derivative / edge.**

We first examine the slowest possible merging of Ω , that is, the latest time of achieving equality of i_{min} , i_{max} . The slowest speed of i_{min} can be 0 (when near the boundary, case 1 of (11)). But i_{max} must always decrease at a speed of at least $2/(l-1)$, through a path of total length no more than $(M-m)$. In this case (Fig. 10d), $t_{merge-in} = (M-m) \cdot (l-1)/2$. Thus,

$$t_{merge-in} < (M-m) \cdot l/2. \quad (12)$$

We now examine the fastest possible merging of Ω with Ω_2 . In this case, i_{min} is at the edge of Ω (Fig. 10e) with zero speed. Thus the merging speed is bounded by the speed of Ω_2 , $2/w_2$, through a path of length of at least $(m-H)$. Thus,

$$t_{merge-out} \geq (m-H) \cdot w_2/2. \quad (13)$$

From the assumption of the Lemma we have:

$$\frac{M-m}{m-H} \leq \frac{w_2}{l} \Rightarrow (M-m) \cdot l \leq (m-H) \cdot w_2. \quad (14)$$

Thus, combining (12), (13) and (14): $t_{merge-in} < t_{merge-out}$. \square

C. Invariance of Texture to Structure

This section gives precise conditions for the validity of Property 4. Let $f : \{0, \dots, N - 1\} \rightarrow \mathbb{R}$ admit the condition defined in Lemma 1. Let $g : \{0, \dots, 2N - 1\} \rightarrow \mathbb{R}$ be a concatenation of f and a translated, value-shifted version of f (see an example in Fig. 10f). We show that texture signatures are identical for fine scales, regardless of their underlying structures. Assuming some constants $0 < C_1 < C_2$, we define:

$$g[i] = \begin{cases} f, & 0 \leq i < N \\ f + C_1, & N \leq i < i_0 + N, i_1 + N + 1 \leq i < 2N \\ f + C_2, & i_0 + N \leq i < i_1 + N + 1. \end{cases}$$

We define the regions of identical texture (up to an additive constant) as $\Omega = \{i_0, \dots, i_1\}$, $\tilde{\Omega} = \{i_0 + N, \dots, i_1 + N\}$.

Theorem 1 (Invariance of Texture to Structure). *Let g be as defined above. Then there exists a time $t_{merge-in}$, such that:*

$$\phi_g(i \in \Omega, t \leq t_{merge-in}) = \phi_g(i \in \tilde{\Omega}, t \leq t_{merge-in}). \quad (15)$$

Proof. First, relying on Lemma 1: $t_{merge-in}(g_{i \in \Omega}) < t_{merge-out}(g_{i \in \Omega})$, and $t_{merge-in}(g_{i \in \tilde{\Omega}}) < t_{merge-out}(g_{i \in \tilde{\Omega}})$.

Second, based on (11), the speed of $u(t)$ at pixel i , $\partial_t u_i$, is invariant to translation and to change by an additive constant. Since $m_{i \in \Omega} = m_{i \in \tilde{\Omega}}$, $\mu_{i \in \Omega} = \mu_{i \in \tilde{\Omega}}$, flow dynamics are identical in Ω and in $\tilde{\Omega}$ until the internal merge.

Therefore: $t_{merge-in} \triangleq t_{merge-in}(g_{i \in \Omega}) = t_{merge-in}(g_{i \in \tilde{\Omega}})$. From the definition of ϕ (4) we deduce:

$$\phi_g(i \in \Omega, t \leq t_{merge-in}) = \phi_g(i \in \tilde{\Omega}, t \leq t_{merge-in}). \quad \square$$

D. Invariance to Linear Illumination Change

This section gives precise conditions for the validity of Property 5. We show that signatures of complete structures, as well as their textures, are invariant to a linear change of baseline, up to a scale related to the scale of the structure (width and height), as seen in Fig. 9i. This is as opposed to the fine scales discussed in Theorem 1, as seen in Fig. 9d.

Let $f : \{0, \dots, N - 1\} \rightarrow \mathbb{R}$ admit the condition defined in Lemma 1, and l be as defined there. Let $g : \{0, \dots, 2N - 1\} \rightarrow \mathbb{R}$ be a concatenation of f and a translated f , with a linearly-changing baseline, as depicted in Fig. 10g, as follows:

$$g[i] = \begin{cases} a \cdot i + b + f[i], & 0 \leq i < N \\ a \cdot i + b + f[i - N], & N \leq i < 2N. \end{cases}$$

We define the structure regions as: $\Omega = \{i_0, \dots, i_1\}$, $\tilde{\Omega} = \{i_2, \dots, i_3\}$, and define: $\Delta \triangleq \max_{i \in \Omega} g_i - \min_{i \in \Omega} g_i$, $w \triangleq i_2 - i_1 - 1$, $h \triangleq g[i_2 - 1] - g[i_1 + 1]$, such that: $a = h/w$. For simplicity, we assume (though this can be relaxed) that structures are of equal heights with respect to the linear slope:

$$\begin{aligned} g[i_0] - g[i_0 - 1] &= g[i_1] - g[i_1 + 1] = \\ g[i_2] - g[i_2 - 1] &= g[i_3] - g[i_3 + 1] \triangleq H. \end{aligned} \quad (16)$$

We assume H is large enough and l small enough, so that local patterns merge first, as in Lemma 1. We also require the following condition:

$$\frac{h(w+1)}{2} > \frac{(H+\Delta) \cdot l}{l+1}. \quad (17)$$

Theorem 2 (Invariance to Linear Illumination Change). *Let g be as defined above. Then there exists a time $t_{merge-struct}$, such that:*

$$\phi_g(i \in \Omega, t \leq t_{merge-struct}) = \phi_g(i \in \tilde{\Omega}, t \leq t_{merge-struct}). \quad (18)$$

Proof. Following the same concept as in Theorem 1, patterns in $\Omega, \tilde{\Omega}$ merge first and simultaneously. To determine the next merging event, we explore the dynamics of structures $\Omega, \tilde{\Omega}$ vs. those of the linear baseline. Using (16), we will now analyze only the dynamics of the left structure and its neighborhood, as both structures behave the same.

For two neighboring pixels $i, j = i + 1$, which are of different regions and of nonzero speed, the merge time according to (11) is $\frac{|g[i] - g[j]|}{2/w_i + 2/w_j}$. Thus the slowest merging of the left structure with the baseline occurs when $\partial_t u_i \neq 0$ only for $i = \{i_0, \dots, i_1 + 1\}$, after its internal evolution made u_{i_1} increase by Δ . Then, an upper bound on the merging time of the left

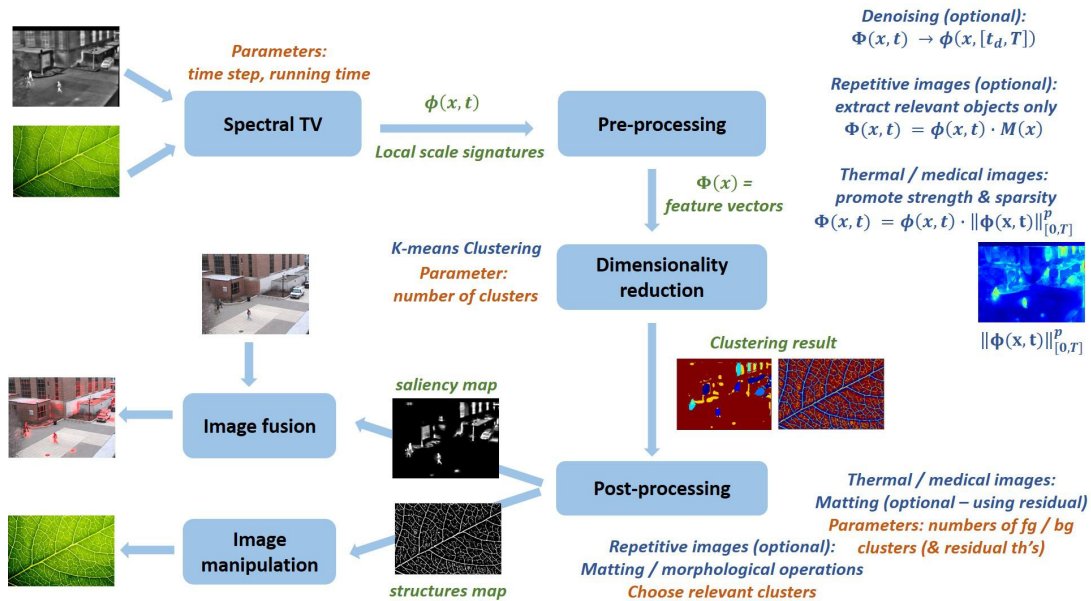


Figure 11: Image manipulation and image fusion using spectral TV local scale signatures: algorithm flowchart.

structure is: $t_{merge-struct} \leq t_{merge}(g_{\{i_1, i_1+1\}}) = \frac{(H+\Delta) \cdot l}{2(l+1)}$. We now calculate the merging time of the linear baseline, $t_{merge-line}$. The preceding evolution is a series of merging events t_k of regions h_k , gradually growing with speeds $v_k = 2/k$ by one pixel at a time, such that $\sum_{k=1}^w h_k = h/w$. Thus: $t_{merge-line} = \sum_{k=1}^w \frac{h_k}{v_k} = \frac{h(1+w)}{4}$. Given (17), $t_{merge-struct} < t_{merge-line}$. Therefore, TV flow dynamics are identical for $\Omega, \tilde{\Omega}$ until $t_{merge-struct}$. Relying on (4), we deduce that:

$$\phi_g(i \in \Omega, t \leq t_{merge-struct}) = \phi_g(i \in \tilde{\Omega}, t \leq t_{merge-struct}). \quad \square$$

VI. ALGORITHM AND APPLICATIONS

The properties of spectral TV signatures can allow grouping together objects with shared features. To facilitate this, we perform dimensionality reduction by clustering of signatures. This allows to partition the image into a set of distinct groups with common spectral TV responses, which can serve for isolating and differentiating salient objects.

A. A Unified Framework

We present a generic unified framework with various applications for images of different modalities (Fig. 11, Algorithm 1). We first decompose the image into its TV elements, using the TV-transform¹ of (4), calculating up to the maximal scale of relevant image structures, T . We use the gray-level version of the input image (in this work, color information is not used). We then perform application-dependent preprocessing on the acquired signatures, denoted $\phi_f(x)$, generating more relevant and enhanced descriptors, denoted $\Phi_f(x)$, to be used as the feature vectors for clustering. A basic dimensionality reduction is performed, using K-means clustering [39]. Last, application-dependent postprocessing of clusters is applied when needed.

¹See spectral TV code in <http://guygilboa.eew.technion.ac.il/code/code1/>

Spectral TV Feature Denoising: Since most image noise appears in the first spectral TV bands, an optional denoising step is inherited within the spectral TV scheme, simply by omitting some of the first spectral TV bands.

Denoting the minimal preserved scale (determined by the expected noise variance) as t_d , and the maximal scale calculated for the transform as T , we define denoising as:

$$\Phi_f(x) = \phi_f(x, [t_d, T]). \quad (19)$$

B. Synthetic Images: Object Differentiation

We first show how disk-type objects are well differentiated in this framework. As we want to use simple unsupervised K-means, we first perform a rough foreground / background separation (as the background may contain several clusters). This is done by exploiting the *negativity* of dominant peaks of background signatures (for dark background). Figs. 3, 13 show clustering results of basic structures. Fig. 12 shows how clustering using our method outperforms clustering using other well-known image decompositions and descriptors [17], [18], [19], [16]. These examples illustrate the unique advantages of the proposed approach. Signatures are very similar for same-object pixels and very distinct compared to pixels of other objects.

C. Application I: Image Manipulation

We extract a map of salient objects of desired sizes or structures for image manipulation: enhancement, attenuation or coloring of certain structures. We can either explore clusters of manually predefined pixels, or choose interesting structures after clustering. Preprocessing may include denoising (Eq. 19), or selecting regions of interest using a map $M(x)$:

$$\Phi_f(x) = \phi_f(x, t) \cdot M(x). \quad (20)$$

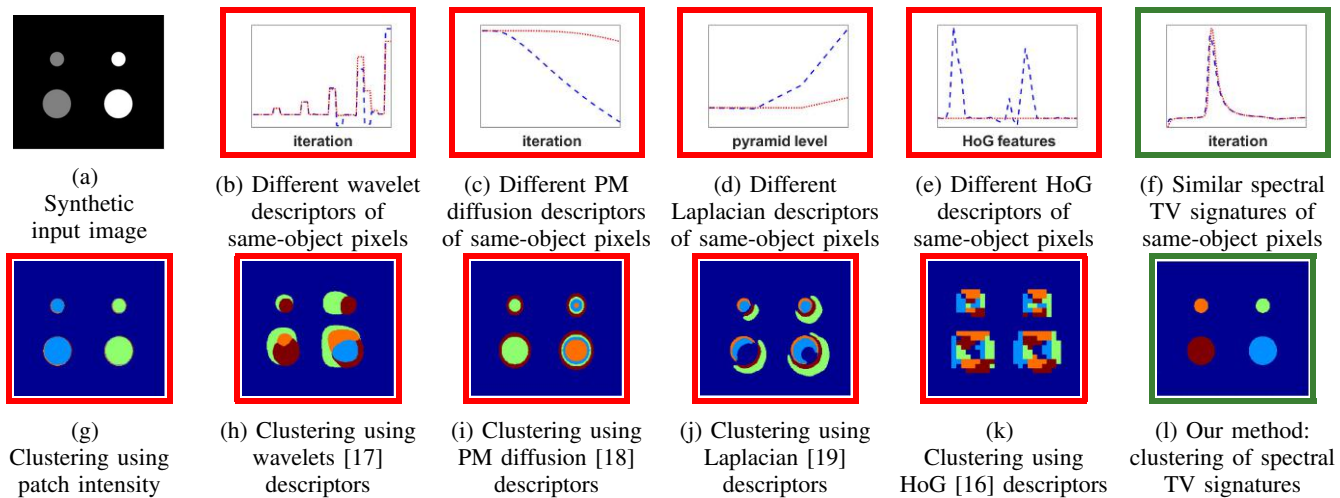


Figure 12: Synthetic analysis and comparison of spectral TV signatures to other descriptors, based on wavelets [17], Perona-Malik (PM) diffusion [18], Laplacian pyramid [19] and Histogram of oriented gradients (HoG) [16]. Clustering (bottom row) succeeds only when relying on spectral TV signatures, since they are the only descriptors (top row), which guarantee similarity for same-object pixels but distinctness for different-object pixels.

Data: Input image: natural with repetitive structures; thermal; medical; or synthetic.

Result: Image manipulation, image fusion or object differentiation.

Spectral TV transform

Preprocessing:

case synthetic image do

Foreground / background separation, Sec. VI-B

case repetitive image do

Spectral TV Denoising, Eq. 19 (optional)

Selecting regions of interest, Eq. 20 (optional)

case thermal or medical image do

Spectral TV Denoising, Eq. 19 (optional)

Signature enhancement, Eq. 21

Dimensionality reduction: K-means clustering

Postprocessing:

case repetitive image do

Matting or morphological operations, Sec. VI-C (optional)

case thermal or medical image do

Matting, Sec. VI-D (optional - using residual, Eq. 6)

Application:

case repetitive image do

Image manipulation

case thermal or medical image do

Image fusion, Appendix B

Algorithm 1: A unified framework for various modalities and applications using spectral TV local scale signatures.

Postprocessing may require image matting [40] or morphological operations of relevant clusters.

D. Application II: Image Fusion

We extract a saliency map from thermal or medical (MRI-T2) images to be fused into a corresponding different-modality

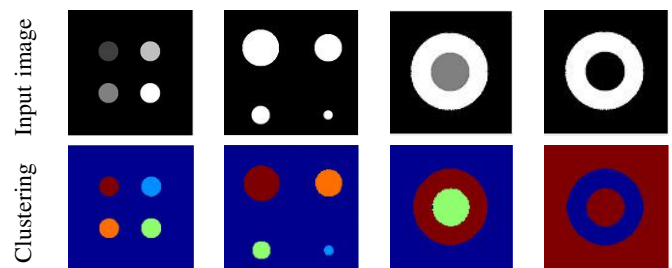


Figure 13: Synthetic images: object differentiation results.

image of the same scene (visible or medical, respectively). This relies on the high contrast of salient objects in these modalities (hot objects, or tumors or abnormal organ structures, respectively). To improve clustering, preprocessing includes signature enhancement - "stretching" each signature according to its L_p norm (usually L_1 norm) over scales:

$$\Phi_f(x) = \phi(x, t) \cdot \|\phi(x, t)\|_{[0, T]}^p. \quad (21)$$

Signatures of salient objects "stretch" more, thus promoting their strength and sparsity (see Fig. 8e vs. Fig. 8f). Denoising (19) is also optional. We cluster these enhanced, possibly denoised signatures, $\Phi_f(x)$. Using more clusters allows capturing smaller or narrower objects, e.g. people and lampposts. Finally, our postprocessing includes image matting [40] to generate the saliency map. The user chooses parameters K_1, K_2 , where the K_1 strongest clusters (in the sense of centroid intensity) form an initial foreground map; and the K_2 weakest clusters form an initial background map. Matting then classifies all other pixels as foreground / background, and the resulting foreground is the saliency map.

The relevance of highly contrasted but large regions is often low, depending on the application. In addition, *weak* signatures may nevertheless indicate *relevant* objects, which have not responded yet within a limited running time, or objects near

image boundaries. We can handle both issues by incorporating large-scale image structures into the postprocessing matting using the image residual $f_r(x)$ (6), requiring to select thresholds for $f_r(x)$. This generates an alternative detailed map with large or near-boundary objects, while avoiding long running times (Figs. 22d vs. 22f). See Appendix B for novel fusion visualization methods.

VII. EXPERIMENTAL RESULTS

We show experimental results for various image modalities and applications, such as image fusion, image segmentation / edge detection and size differentiation, achieving comparable or superior results compared to other techniques².

A. Image Manipulation

Fig. 14 demonstrates salient repetitive structure extraction. A comprehensive comparison clearly shows how our method outperforms other state-of-the-art methods. Our method allows to extract fine salient structures in challenging images: leaf veins with varying-illumination background (b), or thin stripes of a game-board, of the same color as other objects (j). This is thanks to the properties discussed earlier of invariance to rotation, translation and linear illumination. Conversely, other methods [10], [14], [12] fail to detect well such fine structures (c,d,e,k,l,m). Other methods rely on less stable features, which are sensitive to illumination changes (f,n) or to rotations [16] (g,o). In the case of learning-based methods, features are trained for *semantic* segmentation [24] (h,p).

Fig. 15 demonstrates how basic K-means clustering (with $K = 2$) using spectral TV signatures (c) shows highly meaningful clustering compared to the same procedure, based on other well-known image decompositions or descriptors (d,e,f) [17], [18], [19]. Fig. 15b shows an application of signature-based stripe extraction for image manipulation.

The size sensitivity property (Property 1) allows to differentiate objects of similar colors by size. Additional image matting or morphological operations may be used to reconstruct the fine original boundaries from round-shaped clusters. Image manipulation examples based on image segmentation / edge detection or size differentiation are given in Figs. 2, 16.

B. Image Fusion

We give a comprehensive visual and quantitative comparison for multiple state-of-the-art methods, using four established metrics over multiple image datasets. First, we compare the two variants of our method for nine thermal-grayscale image pairs (Table I, Figs. 17 to 20, supplementary). Second, we compare our feature injection method for four medical image pairs (Table II, Fig. 21, supplementary). Third, we compare our temperature gradient coloring method for two thermal-RGB image pairs (Table IV, Figs. 2, 22, supplementary). We also show that our method offers a statistically significant

improvement over other fusion methods. All p-values of the non-parametric Mann-Whitney U test [51] are smaller than 0.05 (Table III).

Our saliency extraction and fusion method presents improved visual results. It can extract fine salient details (Figs. 17, 22), or salient features from a challenging, nearly piecewise constant image (Fig. 20). It also allows extracting differently detailed versions of saliency maps by incorporating large-scale structures (Figs. 20, 22). Our method also outperforms a well established generic saliency extraction method [44], applied directly onto the thermal image (Fig. 17). Note that as opposed to our work, previous thermal saliency extraction work is usually specifically designed for *human* detection [52]. Our fusion scheme also suggests two novel visualization methods, feature injection and gradient coloring (Appendix B).

Our method also achieves superior or comparable quantitative results in almost all cases using the established MI [53], FMI [54] and Petrovic ($Q_P^{AB/F}$) [55] metrics. The prominent exception for this are the inferior results achieved using the Piella measure [56]. Figs. 17 to 19, 21 to 22 visually demonstrate how this metric often does not reflect the visual advantage of our method: maintaining most of the information from the detailed image, while sharply injecting or enhancing the highly-contrasted details.

The drawbacks of the Piella measure have been pointed out before, for example in [57], [54], [58], [59]. More specifically, Cvejic *et al.* [60] have pointed out the reliance of the Piella measure on a problematic definition of image window saliency, and the great influence of the window size on the results. Moreover, Yang *et al.* [61] have pointed out that the Piella measure does not differentiate regions with conflicting information from ones with redundant information. Therefore, methods that visually succeed in selecting the source of information for fusion in these conflicting regions are actually penalized by the measure.

VIII. CONCLUSION

We design an algorithm to isolate and differentiate objects of different contrasts, sizes and structures, as well as multi-scaled objects. We use the comprehensive scale and space information generated by the spectral TV transform, referred to as spectral TV local scale signatures. Given their high dimensionality and redundancy, we reduce their dimensionality to partition an image into meaningful groups. We prove some useful merits of our local framework: sensitivity to size, local contrast and composition of structures, as well as invariance to rotation, translation, flip and linear illumination change. We also provide conditions for the invariance of texture to structure. This enables to construct a unified generic framework applicable for different image modalities and image processing tasks.

APPENDIX

A. Differentiating Compositing Regions by Distinct Signatures

We give an analytic solution for the behavior of the 1D staircase signal to demonstrate Property 2. We show that within each region, signatures of all pixels are identical and distinct with respect to each other (thus can be easily clustered).

²Some demo code will be available at: <https://etyhait.webgr.technion.ac.il/sample-page/spectral-total-variation-local-scale-signatures-for-image-manipulation-and-fusion/>, <http://guygilboa.eew.technion.ac.il/code/>.

Image	Metric	Ours, feature injection	Ours, gradient coloring	GTF [35]	Ratio of LP Pyramid [42]	Wavelets [45], [46]	Complex Wavelets [47]	Curvelet [49]	SVD [48]	MST [41]
Camp	MI[53]	6.8513	5.5641	2.0508	1.4204	1.6191	1.5023	1.4214	1.5417	1.6521
	$Q_P^{AB/F}$ [55]	0.4947	0.4986	0.3761	0.4156	0.2847	0.4245	0.3739	0.2939	0.4884
	FMI [54]	0.5976	0.5737	0.4389	0.4298	0.4309	0.4213	0.3912	0.3226	0.4509
	Q_{Piella} [56]	0.5041	0.5524	0.5100	0.5903	0.6196	0.6139	0.5942	0.5994	0.6481
Dune	MI[53]	6.5156	6.4463	1.3417	1.3586	1.4885	1.3548	1.2754	1.4326	1.5812
	$Q_P^{AB/F}$ [55]	0.5057	0.5061	0.4288	0.4235	0.2983	0.4037	0.3457	0.3381	0.4669
	FMI [54]	0.6146	0.6127	0.5188	0.4386	0.4402	0.4400	0.4076	0.3877	0.4629
	Q_{Piella} [56]	0.4988	0.4997	0.5408	0.6702	0.6886	0.6821	0.6687	0.6844	0.7015
Trees4917	MI[53]	6.3165	6.2372	1.4902	1.4523	1.7048	1.5408	1.4951	1.6130	1.5512
	$Q_P^{AB/F}$ [55]	0.4944	0.4927	0.3578	0.4280	0.2548	0.4043	0.3504	0.2859	0.4524
	FMI [54]	0.5986	0.5966	0.4633	0.4558	0.4176	0.4246	0.3946	0.3280	0.4432
	Q_{Piella} [56]	0.4783	0.4777	0.4669	0.5729	0.6374	0.6382	0.6321	0.6172	0.6374
Road	MI[53]	6.4267	4.7090	1.7227	1.7523	1.7085	1.6805	1.5956	1.5961	2.1023
	$Q_P^{AB/F}$ [55]	0.5796	0.5393	0.4767	0.4942	0.3169	0.5032	0.4631	0.2991	0.5692
	FMI [54]	0.5808	0.5495	0.4717	0.4266	0.4484	0.4329	0.4045	0.2981	0.4668
	Q_{Piella} [56]	0.6292	0.6397	0.6101	0.6724	0.6501	0.6718	0.6417	0.6052	0.7193
Smoke	MI[53]	8.2907	6.0304	3.3017	2.6955	3.3229	2.8532	2.5338	3.1259	4.5191
	$Q_P^{AB/F}$ [55]	0.6411	0.5682	0.3634	0.4890	0.2542	0.5263	0.4755	0.2492	0.5999
	FMI [54]	0.5917	0.5308	0.3986	0.4393	0.4057	0.4408	0.4095	0.2871	0.4739
	Q_{Piella} [56]	0.6245	0.6112	0.5222	0.5917	0.5788	0.6082	0.5576	0.5388	0.6589
T1	MI[53]	7.2481	3.3823	2.2064	2.5584	2.8301	2.2855	2.1993	2.5539	2.7973
	$Q_P^{AB/F}$ [55]	0.6331	0.3406	0.3479	0.5581	0.3656	0.5715	0.5171	0.3638	0.6355
	FMI [54]	0.5803	0.4378	0.4242	0.4409	0.4774	0.4549	0.4223	0.3325	0.4866
	Q_{Piella} [56]	0.7630	0.5723	0.4815	0.7434	0.7046	0.7389	0.7069	0.6805	0.7445
Trees4906	MI[53]	6.4928	6.4088	2.0650	1.9717	2.2998	2.0662	2.0061	2.1789	2.1328
	$Q_P^{AB/F}$ [55]	0.4771	0.4753	0.3509	0.4550	0.2718	0.4525	0.3988	0.2598	0.4904
	FMI [54]	0.5972	0.5949	0.4658	0.4660	0.4268	0.4442	0.4136	0.3019	0.4584
	Q_{Piella} [56]	0.4714	0.4723	0.4808	0.6036	0.6558	0.6645	0.6582	0.6077	0.6719
Steamboat	MI[53]	6.5162	4.8791	3.8905	1.8089	1.9545	1.7626	1.5745	1.7520	4.7079
	$Q_P^{AB/F}$ [55]	0.4424	0.3374	0.2569	0.4386	0.2593	0.4578	0.4167	0.3538	0.5415
	FMI [54]	0.5708	0.5302	0.3849	0.4388	0.4031	0.4239	0.3952	0.3472	0.4762
	Q_{Piella} [56]	0.3137	0.3137	0.4126	0.5729	0.5820	0.5789	0.5501	0.5891	0.6578
Kayak	MI[53]	6.9402	4.2435	2.8485	2.2995	4.2295	2.5632	2.0692	3.4723	2.1464
	$Q_P^{AB/F}$ [55]	0.4568	0.3731	0.1556	0.4983	0.3244	0.6344	0.5476	0.4702	0.6821
	FMI [54]	0.5381	0.4445	0.3203	0.4269	0.4346	0.4158	0.3684	0.3422	0.4500
	Q_{Piella} [56]	0.5209	0.6007	0.2698	0.5912	0.6235	0.6692	0.5534	0.6603	0.7121

Table I: Comprehensive quantitative evaluation of our fusion application. We compare two variants of our method to seven state-of-the-art fusion methods over nine thermal-grayscale image pairs using four established fusion metrics. Our method achieves superior or comparable results almost always (**best result**, **second best result**), except for the Piella measure which we believe does not reflect the visual advantage of our method.

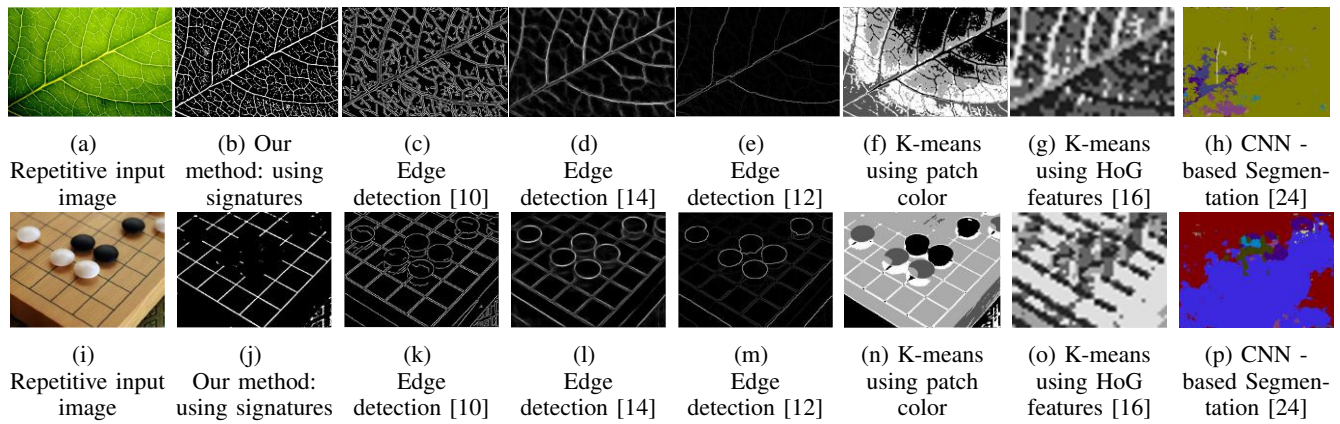


Figure 14: Results and comparisons: salient structure extraction for repetitive images.

Image	Metric	Ours, feature injection	GTF [35]	Ratio of LP Pyramid [42]	Wavelets [45], [46]	Complex Wavelets [47]	Curvelet [49]	SVD [48]	MST [41]
Sarcoma	MI[53]	3.9429	2.7053	2.7470	3.0010	2.5546	2.4985	2.6539	2.8515
	$Q_P^{AB/F}$ [55]	0.5777	0.5529	0.4908	0.4775	0.5218	0.4974	0.4953	0.5766
	FMI [54]	0.7969	0.7816	0.7423	0.7462	0.6395	0.4796	0.7215	0.7231
	Q_{Piella} [56]	0.7769	0.4886	0.8013	0.3447	0.3775	0.3399	0.8153	0.6558
C.T, T2-T1	MI[53]	4.6010	2.6487	3.0589	3.1888	2.7108	2.6509	2.7601	3.1365
	$Q_P^{AB/F}$ [55]	0.5262	0.4575	0.4461	0.3645	0.4833	0.4530	0.3992	0.5207
	FMI [54]	0.7878	0.7459	0.7093	0.6978	0.5809	0.4240	0.6700	0.6934
	Q_{Piella} [56]	0.6362	0.3019	0.7014	0.3045	0.3253	0.3000	0.7125	0.5590
C.T, T2-CT	MI[53]	4.0012	2.3475	2.5217	3.3257	2.5459	2.4296	2.6383	2.7544
	$Q_P^{AB/F}$ [55]	0.4602	0.2718	0.3480	0.3042	0.4498	0.4295	0.3879	0.5162
	FMI [54]	0.7584	0.6972	0.6837	0.7026	0.5148	0.4050	0.6600	0.6696
	Q_{Piella} [56]	0.6844	0.2580	0.6765	0.3473	0.3704	0.3514	0.7507	0.5521
M.B.C	MI[53]	3.8685	2.6531	2.7142	3.7074	2.6783	2.5997	3.0234	2.8625
	$Q_P^{AB/F}$ [55]	0.5414	0.3023	0.4030	0.3650	0.5216	0.5049	0.4537	0.5906
	FMI [54]	0.7977	0.7174	0.7322	0.7540	0.5383	0.4399	0.7034	0.7125
	Q_{Piella} [56]	0.7724	0.4459	0.7569	0.3566	0.4093	0.3962	0.8154	0.6228

Table II: Comprehensive quantitative evaluation of our fusion application. We compare our feature injection method to seven state-of-the-art fusion methods over four medical image pairs (C.T = Cerebral Toxoplasmosis, M.B.C = Metastatic Bronchogenic Carcinoma) using four established fusion metrics. Our method achieves superior or comparable results almost always (**best result**, **second best result**), except for the Piella measure which we believe does not reflect the visual advantage of our method.

Let $f : \{0, \dots, N-1\} \rightarrow \mathbb{R}$ be as depicted in Fig. 10b, and let $0 < i_0 < i_1 < i_2 < i_3 < N-1$. We denote signal regions as $\Omega_1 \triangleq \{i_1, \dots, i_2\}$, $\Omega_2 \triangleq \{i_0, \dots, i_1-1\} \cup \{i_2+1, \dots, i_3\}$, $\Omega_3 \triangleq \{0, \dots, i_0-1\} \cup \{i_3+1, \dots, N-1\}$, and their sizes as $m_1 \triangleq i_2-i_1+1$, $m_2 \triangleq i_1-i_0+i_3-i_2$, $m_3 \triangleq i_0+N-1-i_3$, respectively. Let $f_{i \in \Omega_1} > f_{i \in \Omega_2} > f_{i \in \Omega_3}$. Without loss of generality, we assume that $m_3 > m_1$.

Proposition 1 (Sensitivity to Composition of Structures). *Let f be as defined above. Then:*

$$\phi_f(i \in \Omega_k, t) = \varphi_k(t), \quad k = 1, 2, 3, \quad (22)$$

such that $\varphi_k(t) \neq \varphi_l(t), \forall k \neq l, k, l \in \{1, 2, 3\}$. Note: region signatures are identical $\forall i \in \Omega_k$, even for the disjoint Ω_2 .

Proof. Relying on Section V-A, we analyze the TV flow of f .

1) Phase I: $t \in [0, t_1)$ (Fig. 10b): following (11):

$$u(i, t) = \begin{cases} u(i, 0) + \frac{2t}{m_1} \cdot (-1), & i \in \Omega_1 \\ u(i, 0) + \frac{2t}{m_2} \cdot 0, & i \in \Omega_2 \\ u(i, 0) + \frac{2t}{m_3} \cdot 1, & i \in \Omega_3. \end{cases}$$

$m_3 > m_1 \rightarrow u_t(i \in \Omega_3) < u_t(i \in \Omega_1)$. Thus, at t_1 regions Ω_1, Ω_2 merge to form a new region $\Omega_{1,2}$ of size $m_{1,2} = m_1 + m_2$.

2) Phase II: $t \in [t_1, t_2)$ (Fig. 10c): following (11):

$$u(i, t) = \begin{cases} u(i, t_1) + \frac{2t}{m_{1,2}} \cdot (-1), & i \in \Omega_{1,2} \\ u(i, t_1) + \frac{2t}{m_3} \cdot 1, & i \in \Omega_3. \end{cases}$$

Thus, regions $\Omega_{1,2}, \Omega_3$ merge at $t_2: u(i, t > t_2) = C$.

GTF [35]	Ratio of LP Pyramid [42]	Wavelets [45], [46]	Complex Wavelets [47]	Curvelet [49]	SVD [48]	MST [41]	Spec TV domain fusion [5]
4.1135 e-5	4.1135 e-5	4.1135 e-5	4.1135 e-5	4.1135 e-5	4.1135 e-5	4.1135 e-5	9.99 e-4

Table III: p-values of the Mann-Whitney U test / Wilcoxon rank-sum test [51] for thermal images using the MI metric [53]. All p-values are smaller than 0.05, indicating that our feature injection method offers a statistically significant improvement.

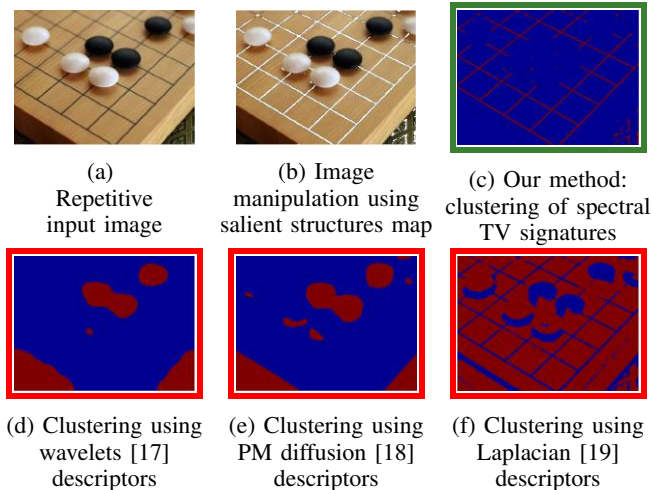


Figure 15: Comparison of K-means clustering into two clusters, using spectral TV signatures vs. other descriptors.

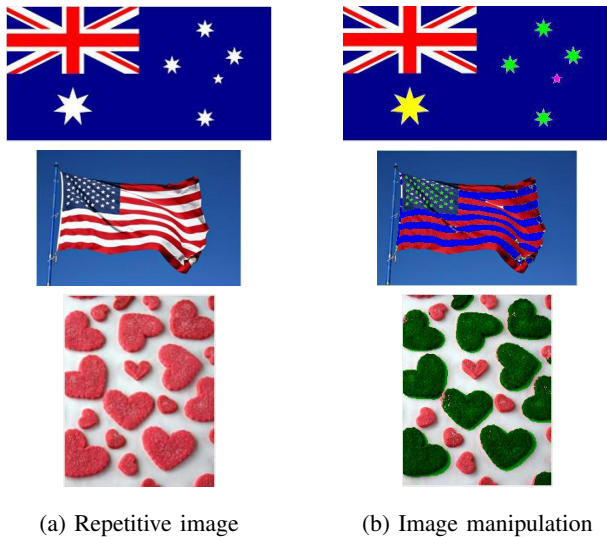


Figure 16: Image manipulation using size differentiation.

Differentiating the TV flow of different regions twice in time:

$$u_{tt}(i, t) = \begin{cases} 2\left(\frac{1}{m_1} - \frac{1}{m_{1,2}}\right)\delta(t - t_1) + \frac{2}{m_{1,2}}\delta(t - t_2), & i \in \Omega_1 \\ \frac{2}{m_{1,2}}\left(-\delta(t - t_1) + \delta(t - t_2)\right), & i \in \Omega_2 \\ -\frac{2}{m_3}\delta(t - t_2), & i \in \Omega_3. \end{cases}$$

From (4) we obtain (for some $A, B, C > 0$):

$$\phi_f(i, t) = \begin{cases} A \cdot \delta(t - t_1) + B \cdot \delta(t - t_2), & i \in \Omega_1 \\ -B \cdot \delta(t - t_1) + B \cdot \delta(t - t_2), & i \in \Omega_2 \\ -C \cdot \delta(t - t_2), & i \in \Omega_3. \end{cases}$$

□

Image	Metric	Ours, gradient coloring	GTF [35]	Wavelets [45], [46]	SVD [48]
Street	MI[53]	7.0235	2.3288	2.1936	1.8259
	$Q_P^{AB/F}$ [55]	0.5291	0.4054	0.3104	0.4053
	FMI [54]	0.5689	0.4261	0.4439	0.3815
	Q_{Piella} [56]	0.5529	0.5503	0.6068	0.6310
Lawn	MI[53]	6.8447	2.5813	3.2628	2.4684
	$Q_P^{AB/F}$ [55]	0.4619	0.4300	0.4323	0.4612
	FMI [54]	0.5492	0.4199	0.4754	0.3428
	Q_{Piella} [56]	0.4714	0.5543	0.6696	0.6582

Table IV: Quantitative evaluation of our fusion application. We compare our temperature gradient coloring method to three state-of-the-art fusion methods over two thermal-RGB image pairs using four established fusion metrics. Our method always achieves superior results (**best result**, **second best result**), except for the Piella measure which we believe does not reflect the visual advantage of our method.

B. Fusion Visualization Methods

Human observers, unlike computer systems, may prefer viewing salient information when fused into a corresponding different-modality image. We suggest two fusion visualization methods. We denote the saliency map as $S(x)$, a corresponding registered image as $V(x)$, and the fused image as $F(x)$. We first suggest injecting salient features directly into the corresponding image (e.g. Fig. 18d):

$$F(x) = \max(V(x), S(x)).$$

This allows introducing information which only appears in $S(x)$ on top of $V(x)$. However, the typically low quality thermal information might overlap the more detailed visible information; and salient white objects will not be visualized as salient. To overcome this, we suggest the temperature gradient coloring method (e.g. Figs. 18c, 21d, 22f). $F(x)$ is a gray-level or RGB replicate of $V(x)$ ($\{R_V(x), G_V(x), B_V(x)\}$), enhanced in locations corresponding to $S(x)$:

$$F(x) = V(x) \cdot g(S(x)),$$

$$\text{or } F(x) = \{R_V(x) \cdot g(S(x)), G_V(x), B_V(x)\},$$

then normalized to avoid clipping. $g(S)$ must be:

- 1) Positive: $\forall x, g(S(x)) > 0$.
- 2) Monotonically increasing: $\forall x, \frac{\partial g(S(x))}{\partial S(x)} > 0$.
- 3) Null for non-salient objects: $g(S(x) = 0) = 1$.

Some useful examples are $g(S) = 1 + S$ and $g(S) = \exp(S)$. Advantages of this method are: avoiding overlapping detailed information with low-quality one; visualizing the gradient of temperatures (e.g., the hotter the object - the redder it appears); and handling salient white objects. Conversely, details which appear only in the saliency map appear weaker in the fusion.

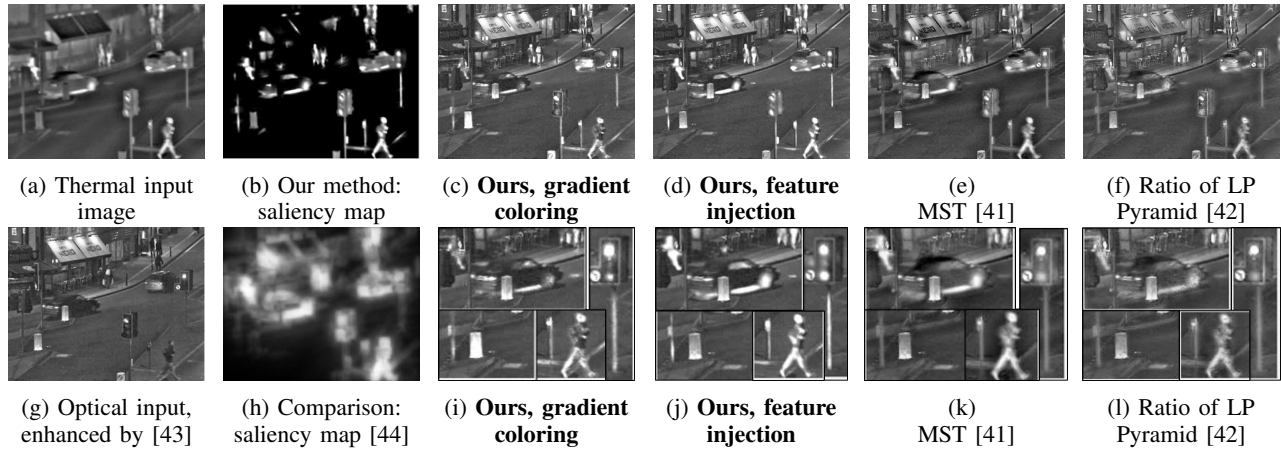


Figure 17: Visual demonstration and comparison for thermal / grayscale fusion (Road image).

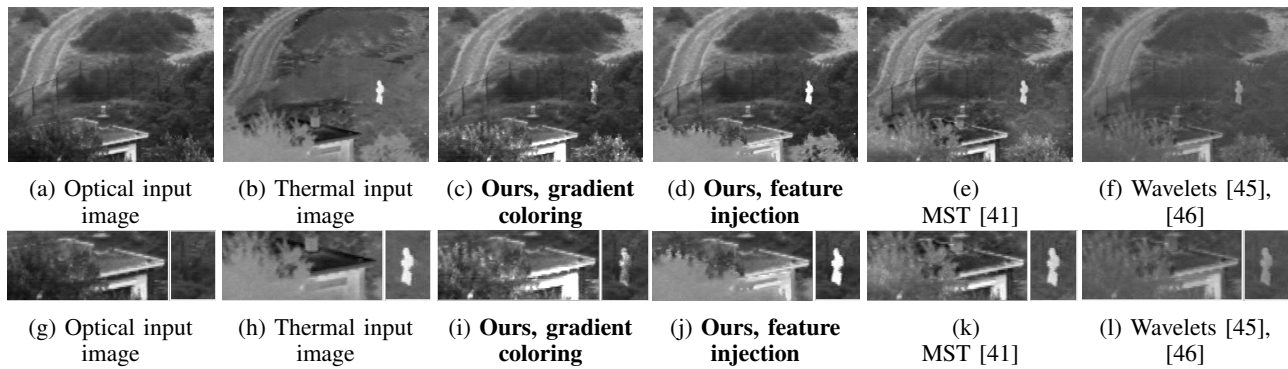


Figure 18: Visual demonstration and comparison for thermal / grayscale fusion (Camp image).

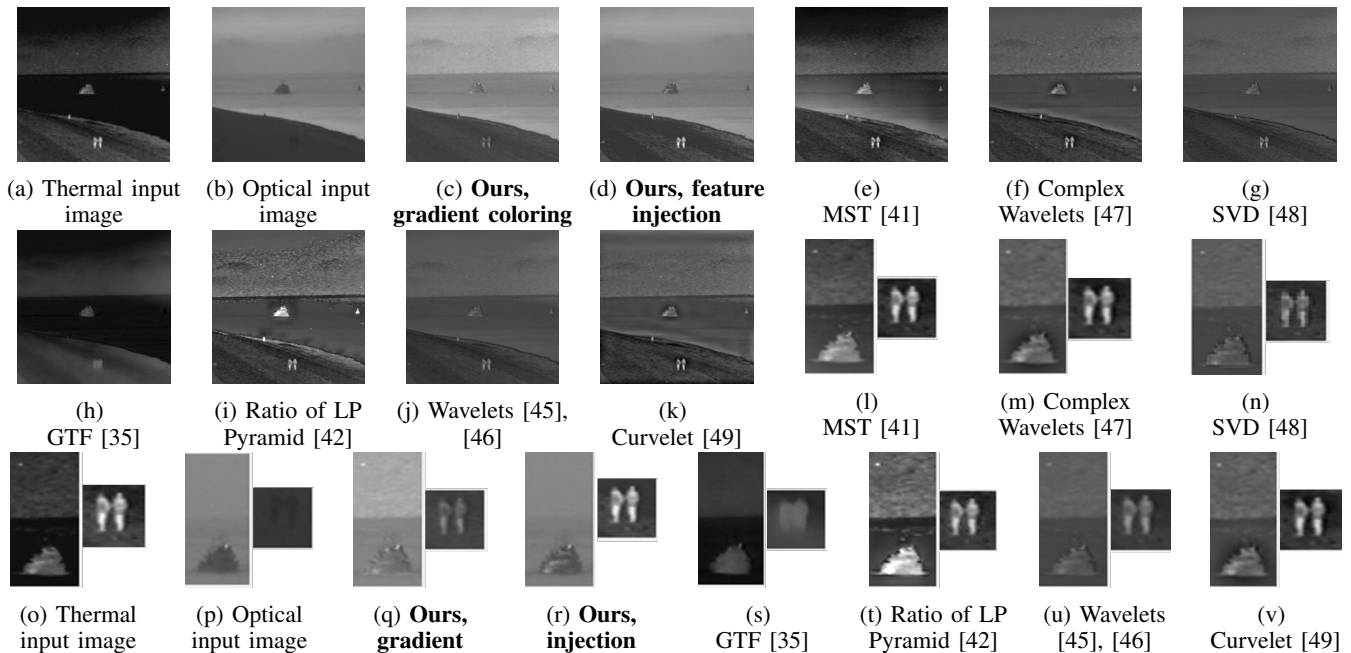


Figure 19: Visual demonstration and comparison for thermal / grayscale fusion (Kayak image).

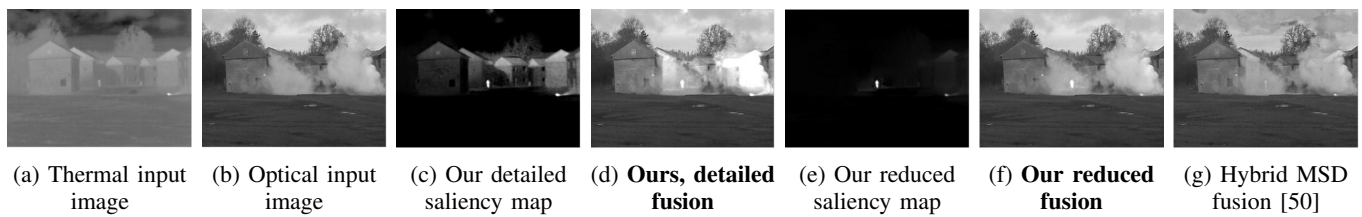


Figure 20: Visual demonstration and comparison for thermal / grayscale fusion (Smoke image).

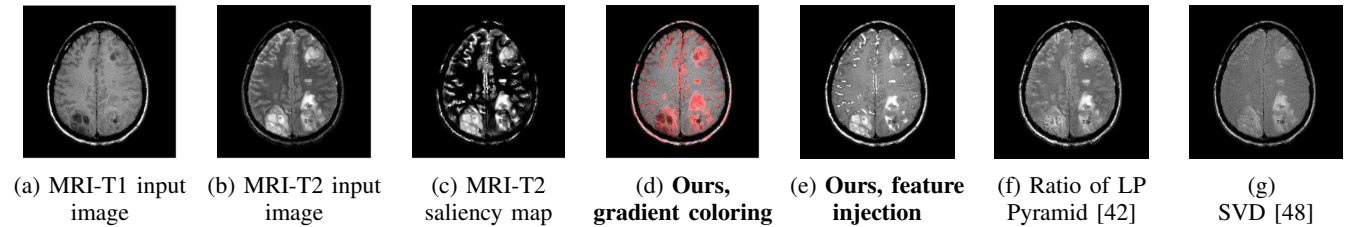


Figure 21: Visual demonstration and comparison for MRI-T2/T1 fusion (Sarcoma image).

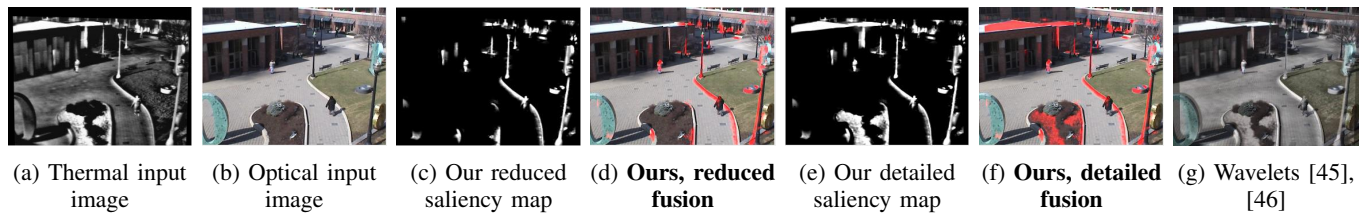


Figure 22: Visual demonstration and comparison for thermal / RGB fusion (Lawn image).

ACKNOWLEDGMENT

Some of the images in this paper were taken from the OTCBVS Benchmark Dataset Collection³, TNO image fusion dataset⁴ and the whole brain Atlas dataset (Harvard)⁵.

REFERENCES

- [1] M. Burger, G. Gilboa, M. Moeller, L. Eckardt, and D. Cremers, "Spectral decompositions using one-homogeneous functionals," *SIAM Journal on Imaging Sciences*, vol. 9, no. 3, pp. 1374–1408, 2016.
- [2] M. Benning, M. Möller, R. Z. Nosssek, M. Burger, D. Cremers, G. Gilboa, and C.-B. Schönlieb, "Nonlinear spectral image fusion," in *International Conference on Scale Space and Variational Methods in Computer Vision*. Springer, 2017, pp. 41–53.
- [3] D. Horesh and G. Gilboa, "Separation surfaces in the spectral tv domain for texture decomposition," *IEEE Transactions on Image Processing*, vol. 25, no. 9, pp. 4260–4270, 2016.
- [4] L. Zeune, S. A. van Gils, L. W. Terstappen, and C. Brune, "Combining contrast invariant l1 data fidelities with nonlinear spectral image decomposition," in *International Conference on Scale Space and Variational Methods in Computer Vision*. Springer, 2017, pp. 80–93.
- [5] H. Lu *et al.*, "Multisensor image fusion and enhancement in spectral total variation domain," *IEEE Transactions on Multimedia*, 2017.
- [6] T. Brox and J. Weickert, "A tv flow based local scale measure for texture discrimination," in *European Conference on Computer Vision*. Springer, 2004, pp. 578–590.
- [7] D. M. Strong, J.-F. Aujol, and T. F. Chan, "Scale recognition, regularization parameter selection, and meyer's g norm in total variation regularization," *Multiscale Modeling & Simulation*, vol. 5, no. 1, pp. 273–303, 2006.
- [8] D. G. Lowe, "Object recognition from local scale-invariant features," in *Computer vision, 1999. The proceedings of the seventh IEEE international conference on*, vol. 2. Ieee, 1999, pp. 1150–1157.
- [9] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool, "Speeded-up robust features (surf)," *Computer vision and image understanding*, vol. 110, no. 3, pp. 346–359, 2008.
- [10] J. Canny, "A computational approach to edge detection," *IEEE Trans. on pattern analysis and machine intelligence*, no. 6, pp. 679–698, 1986.
- [11] D. Comaniciu and P. Meer, "Mean shift: A robust approach toward feature space analysis," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 24, no. 5, pp. 603–619, 2002.
- [12] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik, "Contour detection and hierarchical image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 5, pp. 898–916, 2011.
- [13] D. R. Martin, C. C. Fowlkes, and J. Malik, "Learning to detect natural image boundaries using local brightness, color, and texture cues," *IEEE transactions on pattern analysis and machine intelligence*, vol. 26, no. 5, pp. 530–549, 2004.
- [14] P. Dollár and C. L. Zitnick, "Fast edge detection using structured forests," *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 8, pp. 1558–1570, 2015.
- [15] T. Ojala, M. Pietikäinen, and D. Harwood, "A comparative study of texture measures with classification based on featured distributions," *Pattern recognition*, vol. 29, no. 1, pp. 51–59, 1996.
- [16] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Computer Vision and Pattern Recognition, 2005. CVPR 2005*, vol. 1. IEEE, 2005, pp. 886–893.
- [17] M. Antonini, M. Barlaud, P. Mathieu, and I. Daubechies, "Image coding using wavelet transform," *IEEE Transactions on image processing*, vol. 1, no. 2, pp. 205–220, 1992.
- [18] P. Perona and J. Malik, "Scale-space and edge detection using anisotropic diffusion," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 12, no. 7, pp. 629–639, 1990.
- [19] E. H. Adelson, C. H. Anderson, J. R. Bergen, P. J. Burt, and J. M. Ogden, "Pyramid methods in image processing," *RCA engineer*, vol. 29, no. 6, pp. 33–41, 1984.
- [20] A. K. Jain and F. Farrokhnia, "Unsupervised texture segmentation using gabor filters," *Pattern recognition*, vol. 24, no. 12, pp. 1167–1186, 1991.
- [21] C. Sagiv, N. A. Sochen, and Y. Y. Zeevi, "Integrated active contours for texture segmentation," *IEEE transactions on image processing*, vol. 15, no. 6, pp. 1633–1646, 2006.

³<http://vcipl-okstate.org/pbvs/bench/>

⁴<https://figshare.com/articles/TNO-Image-Fusion-Dataset/1008029>

⁵<http://www.med.harvard.edu/aanlib/>

- [22] S. Xie and Z. Tu, "Holistically-nested edge detection," in *Proc. IEEE international conference on computer vision*, 2015, pp. 1395–1403.
- [23] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *arXiv preprint arXiv:1606.00915*, 2016.
- [24] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 12, pp. 2481–2495, Dec 2017.
- [25] J. Bruna and S. Mallat, "Invariant scattering convolution networks," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1872–1886, 2013.
- [26] S. Mallat, "Understanding deep convolutional networks," *Phil. Trans. R. Soc. A*, vol. 374, no. 2065, p. 20150203, 2016.
- [27] L. I. Rudin, S. Osher, and E. Fatemi, "Nonlinear total variation based noise removal algorithms," *Physica D: Nonlinear Phenomena*, vol. 60, no. 1–4, pp. 259–268, 1992.
- [28] A. Chambolle and P.-L. Lions, "Image recovery via total variation minimization and related problems," *Numerische Mathematik*, vol. 76, no. 2, pp. 167–188, 1997.
- [29] S. Osher, A. Solé, and L. Vese, "Image decomposition and restoration using total variation minimization and the h," *Multiscale Modeling & Simulation*, vol. 1, no. 3, pp. 349–370, 2003.
- [30] G. Gilboa, N. Sochen, and Y. Y. Zeevi, "Estimation of optimal pde-based denoising in the snr sense," *IEEE Transactions on Image Processing*, vol. 15, no. 8, pp. 2269–2280, 2006.
- [31] J.-F. Aujol, G. Gilboa, T. Chan, and S. Osher, "Structure-texture image decomposition: modeling, algorithms, and parameter selection," *International Journal of Computer Vision*, vol. 67, no. 1, pp. 111–136, 2006.
- [32] L. A. Vese and S. J. Osher, "Modeling textures with total variation minimization and oscillating patterns in image processing," *Journal of scientific computing*, vol. 19, no. 1, pp. 553–572, 2003.
- [33] G. Gilboa, N. Sochen, and Y. Y. Zeevi, "Variational denoising of partly textured images by spatially varying constraints," *IEEE Transactions on Image Processing*, vol. 15, no. 8, pp. 2281–2289, 2006.
- [34] M. Kumar and S. Dass, "A total variation-based algorithm for pixel-level image fusion," *IEEE Transactions on Image Processing*, vol. 18, no. 9, pp. 2137–2143, 2009.
- [35] Y. Ma, J. Chen, C. Chen, F. Fan, and J. Ma, "Infrared and visible image fusion using total variation model," *Neurocomputing*, vol. 202, pp. 12–19, 2016.
- [36] F. Andreu, C. Ballester, V. Caselles, J. Mazón *et al.*, "Minimizing total variation flow," *Differential and integral equations*, vol. 14, no. 3, pp. 321–360, 2001.
- [37] G. Gilboa, "A total variation spectral framework for scale and texture analysis," *SIAM journal on Imaging Sciences*, vol. 7, no. 4, pp. 1937–1961, 2014.
- [38] G. Steidl, J. Weickert, T. Brox, P. Mrázek, and M. Welk, "On the equivalence of soft wavelet shrinkage, total variation diffusion, total variation regularization, and sides," *SIAM Journal on Numerical Analysis*, vol. 42, no. 2, pp. 686–713, 2004.
- [39] D. Arthur and S. Vassilvitskii, "k-means++: The advantages of careful seeding," in *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*. Society for Industrial and Applied Mathematics, 2007, pp. 1027–1035.
- [40] A. Levin, D. Lischinski, and Y. Weiss, "A closed-form solution to natural image matting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 30, no. 2, pp. 228–242, 2008.
- [41] Y. Liu, S. Liu, and Z. Wang, "A general framework for image fusion based on multi-scale transform and sparse representation," *Information Fusion*, vol. 24, pp. 147–164, 2015.
- [42] A. Toet, "Image fusion by a ratio of low-pass pyramid," *Pattern Recognition Letters*, vol. 9, no. 4, pp. 245–253, 1989.
- [43] Z. Zhou, M. Dong, X. Xie, and Z. Gao, "Fusion of infrared and visible images for night-vision context enhancement," *Applied optics*, vol. 55, no. 23, pp. 6480–6490, 2016.
- [44] S. Goferman, L. Zelnik-Manor, and A. Tal, "Context-aware saliency detection," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 10, pp. 1915–1926, 2012.
- [45] P. Zeeuw, "Wavelet and image fusion," *CWI, Amsterdam, March*, vol. 444, 1998.
- [46] M. Misiti, "Les ondelettes et leurs applications. pdf," 2003.
- [47] J. J. Lewis, R. J. OCallaghan, S. G. Nikolov, D. R. Bull, and N. Canagarajah, "Pixel-and region-based image fusion with complex wavelets," *Information fusion*, vol. 8, no. 2, pp. 119–130, 2007.
- [48] V. Naidu, "Image fusion technique using multi-resolution singular value decomposition," *Defence Science Journal*, vol. 61, no. 5, pp. 479–484, 2011.
- [49] F. Nencini, A. Garzelli, S. Baronti, and L. Alparone, "Remote sensing image fusion using the curvelet transform," *Information fusion*, vol. 8, no. 2, pp. 143–156, 2007.
- [50] Z. Zhou, B. Wang, S. Li, and M. Dong, "Perceptual fusion of infrared and visible images through a hybrid multi-scale decomposition with gaussian and bilateral filters," *Information Fusion*, vol. 30, pp. 15–26, 2016.
- [51] H. B. Mann and D. R. Whitney, "On a test of whether one of two random variables is stochastically larger than the other," *The annals of mathematical statistics*, pp. 50–60, 1947.
- [52] J. W. Davis and V. Sharma, "Robust background-subtraction for person detection in thermal imagery," in *CVPR Workshops*, 2004, p. 128.
- [53] G. Qu, D. Zhang, and P. Yan, "Information measure for performance of image fusion," *Electronics letters*, vol. 38, no. 7, pp. 313–315, 2002.
- [54] M. B. A. Haghighat, A. Aghagolzadeh, and H. Seyedarabi, "A non-reference image fusion metric based on mutual information of image features," *Computers & Electrical Engineering*, vol. 37, no. 5, pp. 744–756, 2011.
- [55] C. Xydeas, , and V. Petrovic, "Objective image fusion performance measure," *Electronics letters*, vol. 36, no. 4, pp. 308–309, 2000.
- [56] G. Piella and H. Heijmans, "A new quality metric for image fusion," in *Image Processing, 2003. ICIP 2003. Proceedings. 2003 International Conference on*, vol. 3. IEEE, 2003, pp. III–173.
- [57] N. Cvejic, T. Seppanen, and S. J. Godsill, "A nonreference image fusion metric based on the regional importance measure," *IEEE Journal of Selected Topics in Signal Processing*, vol. 3, no. 2, pp. 212–221, 2009.
- [58] P.-w. Wang and B. Liu, "A novel image fusion metric based on multi-scale analysis," in *Signal Processing, 2008. ICSP 2008. 9th International Conference on*. IEEE, 2008, pp. 965–968.
- [59] M. Haghighat and M. A. Razian, "Fast-fmi: non-reference image fusion metric," in *Application of Information and Communication Technologies (AICT), 2014 IEEE 8th International Conference on*. IEEE, 2014, pp. 1–3.
- [60] N. Cvejic, A. Loza, D. Bull, and N. Canagarajah, "A similarity metric for assessment of image fusion algorithms," *International journal of signal processing*, vol. 2, no. 3, pp. 178–182, 2005.
- [61] C. Yang, J.-Q. Zhang, X.-R. Wang, and X. Liu, "A novel similarity based quality metric for image fusion," *Information Fusion*, vol. 9, no. 2, pp. 156–160, 2008.

Ester Hait is pursuing her Ph.D. in Electrical Engineering in the Technion - Israel Institute of Technology, where she received her B.Sc. (Cum Laude) and M.Sc. in Electrical Engineering in 2014 and 2016, respectively. Her research interests include image processing and computer vision.

Guy Gilboa received his Ph.D. from the Electrical Engineering Department, Technion - Israel Institute of Technology in 2004. He was a postdoctoral fellow with UCLA and had various development and research roles with Microsoft and Philips Healthcare. Since 2013 he is a faculty member with the Electrical Engineering Department, Technion - Israel Institute of Technology. He has authored some highly cited papers on topics such as image sharpening and denoising, nonlocal operators theory, and texture analysis. He received several prizes, including the Eshkol Prize by the Israeli Ministry of Science, the Vatav Scholarship, and the Gutwirth Prize. He serves at the editorial boards of the journals IEEE SPL, JMIV and CVIU.