

Stochastic Image Models from SIFT-Like Descriptors*

A. Desolneux[†] and A. Leclaire[†]

Abstract. Extraction of local features constitutes a first step of many algorithms used in computer vision. The choice of keypoints and local features is often driven by the optimization of a performance criterion on a given computer vision task, which sometimes makes the extracted content difficult to apprehend. In this paper we propose to examine the content of local image descriptors from a reconstruction perspective. For that, relying on the keypoints and descriptors provided by the scale-invariant feature transform (SIFT), we propose two stochastic models for exploring the set of images that can be obtained from given SIFT descriptors. The two models are both defined as solutions of generalized Poisson problems that combine gradient information at different scales. The first model consists in sampling an orientation field according to a maximum entropy distribution constrained by local histograms of gradient orientations (at scale 0). The second model consists in simple resampling of the local histogram of gradient orientations at multiple scales. We show that both of these models admit convolutive expressions which allow us to compute the model statistics (e.g., the mean, the variance). Also, in the experimental section, we show that these models are able to recover many image structures, while not requiring any external database. Finally, we compare several other choices of points of interest in terms of quality of reconstruction, which confirms the optimality of the SIFT keypoints over simpler alternatives.

Key words. image synthesis, random image model, reconstruction from features, SIFT, Poisson editing, maximum entropy distributions, exponential models

AMS subject classifications. 62M40, 65D18, 68U10, 94A08

DOI. 10.1137/18M116592X

1. Introduction. A fundamental problem of vision consists in extracting a minimal representation that is sufficient for a human to apprehend the semantic content of an image. Marr and Hildreth [39, 38] proposed a *raw primal sketch* image representation based on the zero-crossings of the Laplacian computed at different scales, which extract spatial positions corresponding to edges, blobs, and terminations. Since this pioneering work, many authors proposed extracting different points of interest (keypoints) or local descriptors (features) based on several differential operators, while being invariant to given image transformations. Extracting keypoints and local features in images is indeed a fundamental step for many imaging tasks [21], like image recognition [63, 32, 9, 10, 57], image matching and rectification [32, 60, 31], object detection and tracking [8, 58, 66, 52], video stabilization [6, 65], and image classification [28, 68, 27]. In this paper, we propose to discuss the role of such keypoints and descriptors from a reconstruction point of view.

In the seminal paper [5], Attneave suggests that the most important points for image

*Received by the editors January 19, 2018; accepted for publication (in revised form) July 23, 2018; published electronically October 16, 2018. A preliminary version of this work was published as a conference paper in [17].

<http://www.siam.org/journals/siims/11-4/M116592.html>

[†]CMLA, ENS Cachan, CNRS, Université Paris-Saclay, 94235 Cachan, France (agnes.desolneux@cmla.ens-cachan.fr, arthur.leclaire@cmla.ens-cachan.fr).

perception are the ones of maximum curvature. Since then, many techniques have emerged to single out keypoints and build local descriptors around them. Depending on the applicative context, one should use descriptors that are invariant with respect to specific geometric transformations¹ (e.g., image recognition generally needs invariance to homography and illumination change). Here we will only mention a few famous local descriptors, and we refer the reader to [42, 59, 44, 31] for a more comprehensive survey.

Harris and Stephens proposed a combined corner and edge detector based on the determinant and trace of the structure tensor of the image [23]. A multiscale variant based on a normalized Laplacian of Gaussian (LoG) scale-space, coined Harris–Laplace, was proposed by Milokajczyk and Schmid [41]. The same authors also proposed in [41] the Harris-affine point detector, which extends the previous one with a normalization step in order to get invariance to affine transformations. Tuytelaars and Mikolajczyk proposed in [59] two region detectors, both starting from anchor points (e.g., Harris points); then the first one selects a region within detected edges around the anchor, and the second one extracts a region by analyzing intensity profiles on rays emanating from the anchor. Rosten and Drummond introduced in [54] the “features from accelerated segment test” (FAST), which is a corner detector accelerated by a machine learning technique. This approach was made faster by Mair et al. [36] using optimal decision trees, thus obtaining an “adaptive and generic accelerated segment test” (AGAST). Musé et al. [47] proposed extracting shapes from the image level lines, and processing them in order to get an affine invariant representation.

In parallel to this research on keypoints, many techniques have been proposed for invariant local descriptions of images. An early descriptor is given by the local binary patterns (LBPs) defined by Ojala, Pietikäinen, and Mäenpää [50], which extract signs of differences of image values on pixels located on a circular neighborhood of a keypoint. The LBPs were originally designed for texture description but can also be used for face detection [1]. In [32], Lowe introduced the scale-invariant feature transform (SIFT), which first extracts the keypoints as local extrema of the “difference of Gaussian” (DoG) approximation of the LoG, and next computes around each keypoint a local descriptor based on normalized histograms of gradient direction (HOG); see the details in section 2. Notice that similar HOG descriptors computed on a dense grid were actually used in [14] for person detection; one reference implementation of the HOG descriptors is given in [22]. A fully affine-invariant extension of SIFT, named ASIFT, was proposed by Morel and Yu [44] and consists in applying the SIFT method with the image transformed with several simulated affine maps. The SURF method (speeded-up robust features) proposed by Bay et al. [7] is closely related in construction to the SIFT method, but allows for a faster implementation. At a higher semantic level, local image behavior can also be represented by visual words [58, 11] which are obtained as cluster points in a feature space. Later, some authors proposed describing a patch using local binary descriptors (LBDs), which extract the signs of differences between Gaussian measurements taken at different locations. Using different ways of selecting these locations leads to the methods BRISK [29] (binary robust invariant scalable keypoints) or FREAK [2] (fast retina keypoint). All of these descriptors have quite different invariance properties (evaluated in either a theoretical or an experimental framework).

¹The translation invariance is generally always required, and often trivial.

Long before the design of these image descriptors, the question of a minimal representation of an image was thoroughly studied, mainly for compression purposes. Through the concept of a *raw primal sketch*, Marr [38] suggested that the human visual system processes images by retaining essentially the lines of zero-crossing of the Laplacian at several scales. This led to the conjecture that an image is uniquely defined by these zero-crossing lines, a conjecture that was later made precise by Mallat and Zhong [37] using wavelet modulus maxima. Both of these conjectures were proved wrong by Meyer [40], but still algorithms for approximate reconstruction were proposed by Hummel and Moniot [24] for zero-crossings and by Mallat and Zhong [37] for the case of wavelet modulus maxima. In addition, unique characterization can be shown to be true under some additional hypotheses [12, 13, 55, 4, 3].

From a more practical point of view, several authors have raised the question of inversion of a feature-based representation. For example, Elder and Zucker [20] proposed an algorithm for image reconstruction from detected contours based on the heat diffusion. Nielsen and Lillholm [49] considered the problem of variational reconstruction from linear measurements; in addition to the minimum variance reconstruction (given by the pseudoinverse of the measurements matrix), they proposed two variational reconstructions based on either the entropy (of the image seen as a probability distribution on its domain) or the H^1 norm. Interestingly, they discussed the problem of extracting a subset of linear measurements leading to the best reconstruction and empirically compared three different strategies for that purpose.

Motivated by privacy issues (since the descriptors may be transmitted on an unsecured network), Weinzaepfel, Jégou, and Pérez [64] addressed image reconstruction from the output of a SIFT transform adapted with elliptic keypoints. One important difference with previous works is that this method exploits a database of image patches: for each keypoint, a patch with a similar description is sought in the database, and all the patches are stitched together with Poisson image editing [51]. Vondrick et al. [62] addressed reconstruction from dense HOGs by relying on a paired dictionary representation of HOGs and patches. Also, d'Angelo et al. [15] addressed reconstruction from local binary descriptors by relying on primal-dual optimization techniques; in contrast with [64, 62], this method does not need any external information. Kato and Harada [26] formulated reconstruction from a bag of visual words as a problem of quadratic assignment. Finally, Juefei-Xu and Savvides [25] proposed inverting the LBP representation with an approach based on paired dictionary learning with an ℓ^0 constraint.

More recently, the success of deep convolutional neural networks in image classification [27, 67] has urged the need for inverting the corresponding representations in order to intuitively understand the kind of information that is extracted at each layer. Even if they do not formulate it as an inverting procedure, Zeiler and Fergus [67] proposed building a deconvolution network that allows one to visualize in image space the stimuli that excite one response at a particular layer of the neural network. Given an image u , Mahendran and Vedaldi [34, 35] proposed searching for a preimage of an image representation $\varphi(u)$ by minimizing a functional containing a loss term related to the representation φ and a regularizing term (in particular the H^1 norm). Even if the regularizer is convex, the transformation φ is in general highly nonlinear, so that the resulting optimization problem is not convex; thus the output of the inversion may depend on the parameters and initializations of the chosen optimization procedure. On the other hand, Dosovitskiy and Brox [19] suggested learning an approximate left

inverse of the representation (i.e., a mapping φ_L^{-1} such that $\varphi_L^{-1}(\varphi(u)) \approx u$ for every u) in the form of an up-convolutional network. These methods are generic in the sense that they can be applied to any image representation that can be approximated by the output of a convolutional neural network; in particular, the authors of [19] display inversion results for HOG, SIFT, and AlexNet [27] representations. Notice that the inversion/visualization techniques of [67, 19] exploit an external database, while that of [34, 35] does not.

Instead of building a uniquely defined inversion technique (using regularization), another way to perform reconstruction from the image representation φ is to sample from a stochastic model that explores the set of preimages of $\varphi(u)$. This is particularly relevant if one uses an image representation that is not invertible: for example, the SIFT cells of an image may not cover its whole domain, and thus many images could have the same SIFT descriptors. Besides, the HOG descriptors are inherently of a statistical nature: each HOG extracts the distribution of gradient orientations in one small area. Thus they only provide a locally pooled information and thus do not precisely constrain each gradient value. For this reason, the inversion by direct (regularized) optimization proposed in [34, 35] is not adapted to the usual SIFT representation (sometimes called sparse SIFT, as opposed to SIFT descriptors computed on a dense grid).

One way to address this problem is to sample from a maximum entropy model that complies with these statistical constraints. Such maximal entropy models were considered by Zhu, Wu, and Mumford in [69, 46] for texture modeling based on responses to an automatically selected subset of filters chosen from a filter bank. This approach has been recently extended by Lu, Zhu, and Wu to responses to a pretrained neural network [33]. Maximum entropy models were also used to question the noise models used in the *a contrario* framework for feature detections in images [18]: in [16], for two types of given detections (cluster of points, or line segments), Desolneux proposed explicit computations of maximal entropy image models that led to the same detections (in average). Let us emphasize that one important difference with previous works is that, more than reconstructing the original image, we aim at exploring the set of images with similar HOG description at the keypoints positions, with the least possible a priori information on what the reconstruction should look like. In contrast, the dependence on an external database in [51, 19] poses a strong prior on the reconstruction.

In the present paper, we propose two stochastic models that comply with statistical features given by a SIFT-like representation. In order to derive explicit computations, we work on a simplified SIFT transform which extracts multiscale HOGs from regions around the (usual) SIFT keypoints. The first model, called MaxEnt, is indeed an instance of a maximum entropy model which complies with local statistical constraints on the gradient orientations (at scale 0, i.e., the image scale). Once the parameters of this model are estimated (using a gradient descent), a target gradient orientation can be sampled, and we recover an image by solving a classical Poisson problem. The second model, called MS-Poisson, consists in first independent sampling of multiscale gradient orientations in all the SIFT cells, and next merging all the pieces by solving a global multiscale Poisson problem. Even if this model does not solve an explicit maximum entropy problem, it allows one to coherently merge information given at several scales. Several experiments show that both models are able to recover large image structures and compare well to the results of [64] while not using any external information. Finally, we discuss the definition of the SIFT keypoints in terms of optimality of

reconstruction, thus raising the following question related to visual information theory: “Can we measure the optimality (at a fixed memory budget) of some image descriptor in terms of reconstruction?”

The paper is organized as follows. In section 2, we briefly recall the main steps of the SIFT method and explain the simplified SIFT descriptors that we use for reconstruction. In section 3, we build and study the maximum entropy model (MaxEnt) used for reconstruction from monoscale HOGs computed in the SIFT subcells. In section 4, we propose the multiscale Poisson model (MS-Poisson) that allows us to comply with multiscale HOGs taken in the SIFT subcells; the corresponding H^1 -regularized multiscale Poisson problem is explicitly solved. In section 5 we display several reconstruction results obtained with both models (applied with simplified SIFT, or also the true SIFT) and study the variability of the reconstruction (in terms of first and second order moments, but also of SIFT keypoints computed on the reconstruction). We also compare our results with other existing reconstruction techniques and apply the reconstruction models on other keypoint sets, thus confirming (from the synthesis perspective) the efficiency of the SIFT method for global image description. Finally, in section 6 we conclude the discussion proposed in this paper and offer some possibilities for future research. A preliminary version of this work was published as a conference paper in [17]. Compared to the conference version, here we explain in more detail the derivation of MaxEnt and MS-Poisson models, providing some more properties of these models and in particular explicit formulae for the first and second order moments of these models. We also propose several new experiments which illustrate the performance and the variability of these models (with qualitative and quantitative evaluation) and question the role of the keypoint definition in the quality of reconstruction. Also, the supplementary material (M116592_01.zip [local/web 116MB]) attached to this paper contain source codes which allow the reader to reproduce the experiments shown in the paper.

2. A brief summary of the SIFT method. In this section we briefly recall the construction of keypoints and local descriptors used in the SIFT method, and we explain the simplified descriptors that will be used later for the reconstruction in the next sections.

2.1. Gaussian scale-space and keypoints. Following [30], we introduce the Gaussian scale-space in a continuous domain. Let $u : \mathbb{R}^2 \rightarrow \mathbb{R}$ be an integrable function. For $\sigma > 0$, we introduce the function $g_\sigma : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by

$$g_\sigma(\mathbf{x}) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{|\mathbf{x}|^2}{2\sigma^2}\right).$$

The Gaussian scale-space associated with u is then defined by the convolution

$$\forall \mathbf{x} \in \mathbb{R}^2, \forall \sigma > 0, \quad L_u(\mathbf{x}, \sigma) = g_\sigma * u(\mathbf{x}) = \int_{\mathbb{R}^2} g_\sigma(\mathbf{y})u(\mathbf{x} - \mathbf{y})d\mathbf{y}.$$

Another way to parameterize the scale-space is to use a time parameter $t = \sigma^2$ and the kernel $k_t = g_{\sqrt{t}}$, which satisfies

$$\frac{\partial}{\partial t}(k_t(\mathbf{x})) = \frac{1}{2}\Delta k_t(\mathbf{x}).$$

In other words, $(\mathbf{x}, t) \mapsto L_u(\mathbf{x}, \sqrt{t})$ is the solution of the heat equation on \mathbb{R}^2 with initial condition u (in particular, it is a C^∞ function on $\mathbb{R}^2 \times (0, \infty)$).

Then we consider the scale-normalized Laplacian of Gaussian (LoG) $\sigma^2 \Delta g_\sigma$. The PDE satisfied by k_t gives, after a change of variables, that

$$\sigma \frac{\partial g_\sigma}{\partial \sigma}(\mathbf{x}) = \sigma^2 \Delta g_\sigma(\mathbf{x}) = \left(\frac{|\mathbf{x}|^2 - 2\sigma^2}{2\pi\sigma^4} \right) \exp\left(-\frac{|\mathbf{x}|^2}{2\sigma^2}\right).$$

The detection of keypoints will be based on the local extrema of the function

$$D_u(\mathbf{x}, \sigma) := \sigma^2 \Delta g_\sigma * u(\mathbf{x}) = \sigma^2 \Delta (g_\sigma * u)(\mathbf{x}).$$

The following proposition, which is recalled without proof, shows that these keypoints are covariant to several image transformations.

Proposition 1 (see [30]). *We have the following invariance properties.*

1. $\forall a \in \mathbb{R}$, $D_{au} = aD_u$.
2. If v is an affine function of \mathbf{x} , then $D_{u+v} = D_u$.
3. If $\mathbf{h} \in \mathbb{R}^2$ and $\tau_{\mathbf{h}}u(\mathbf{x}) = u(\mathbf{x} - \mathbf{h})$ is a translated version of u , then

$$D_{\tau_{\mathbf{h}}u}(\mathbf{x}, \sigma) = D_u(\mathbf{x} - \mathbf{h}, \sigma).$$

4. (Scale invariance.) If $u(\mathbf{x}) = v(s\mathbf{x})$ with $s > 0$, for all $\mathbf{x} \in \mathbb{R}^2$, then

$$D_u(\mathbf{x}, \sigma) = D_v(s\mathbf{x}, s\sigma).$$

The existence of a keypoint (\mathbf{x}, σ) indicates the presence of a blob-like structure at position \mathbf{x} with scale σ . For example, the Gaussian function g_s ($s > 0$) admits a keypoint $(0, s)$ which corresponds to a strict local minimum of D_{g_s} .

The authors of [45] also discussed the effect of several other image transformations on the SIFT keypoints but left aside the factor σ^2 in the definition of D_u .

2.2. SIFT summary. In the paper by Lowe [32], the scale-normalized LoG is approximated by a finite difference of Gaussian (DoG) functions: for a constant scale factor $k > 1$, he considers instead

$$(1) \quad (\mathbf{x}, \sigma) \mapsto (g_{k\sigma} - g_\sigma) * u(\mathbf{x}) \approx (k\sigma - \sigma) \frac{\partial g_\sigma}{\partial \sigma} * u(\mathbf{x}) = (k-1)\sigma^2 \Delta g_\sigma * u(\mathbf{x}).$$

Also, the practical implementation of [32] only works with discretized images, so that the extracted keypoints are actually strict local extrema computed on a discretized scale-space.

Here is a quick summary of the original SIFT method [32]. For technical details we refer the reader to [53]. Here, and in the remainder of the paper, u_0 refers to the original image on which we compute keypoints and local descriptors.

1. Computing SIFT keypoints:
 - (a) Extract local extrema of a discrete version of (1).
 - (b) Refine the positions of the local extrema in position and scale using a quadratic approximation.

- (c) Discard extrema with low contrast (thresholding low values of (1)) and extrema located on edges (thresholding high values of the ratio between Hessian eigenvalues).
2. Computing SIFT local descriptors associated with the keypoint (\mathbf{x}, σ) :
- (a) Compute one or several principal orientations α . For that, in a square of size $9\sigma \times 9\sigma$ centered at \mathbf{x} (and parallel to the image axes), compute a smoothed histogram of orientations of $\nabla g_\sigma * u_0$ and extract its significant local maxima.
- (b) For each detected orientation α , consider a grid of 4×4 square regions around (\mathbf{x}, σ) . These square regions, which we call SIFT subcells, are of size $3\sigma \times 3\sigma$ with one side parallel to α . In each subcell compute the histogram of $\text{Angle}(\nabla g_\sigma * u_0) - \alpha$ quantized on 8 values ($\ell \frac{\pi}{4}, 1 \leq \ell \leq 8$).
- (c) Normalization: the 16 histograms are concatenated to obtain a feature vector $f \in \mathbb{R}^{128}$, which is thresholded and normalized,

$$(2) \quad f_k \leftarrow \min(f_k, 0.2\|f\|_2), \quad f_k = \min\left(255, \left\lfloor 512 \frac{f_k}{\|f\|_2} \right\rfloor\right),$$

and finally quantized to 8-bit integers.

When computing orientation histograms in steps 2(a) and 2(b), each pixel votes with a weight that depends on the value of the gradient norm at scale σ and on its distance to the keypoint center \mathbf{x} . Also in step 2(b), there is a linear splitting of the vote of an angle between the two adjacent quantized angle values.

2.3. Keypoints and descriptors used in our method. In the reconstruction models proposed in this paper, we work with images defined on a rectangle $\Omega \subset \mathbb{Z}^2$ and consider the oriented keypoints extracted by the original SIFT method. However, we will only work with simplified SIFT descriptors in the sense that we extract hard-binned histograms of gradient orientations at several scales. In other words, we do not include the vote weights or the normalization step 2(c).

We thus denote by $(s_j)_{j \in \mathcal{J}}$ the collection of SIFT subcells, $s_j \subset \Omega$ (if a $3\sigma \times 3\sigma$ subcell is not entirely contained in Ω , then we replace it with its intersection with Ω). The SIFT *subcells* must not be confused with the SIFT *cells*: in a SIFT cell, there are 16 SIFT subcells, so that different subcells s_j can correspond to the same keypoint. We will denote by $(\mathbf{x}_j, \sigma_j, \alpha_j)$ the oriented keypoint associated with s_j . For $\mathbf{y} \in \Omega$, we denote by $\mathcal{J}(\mathbf{y}) = \{j \in \mathcal{J} \mid \mathbf{y} \in s_j\}$ the set of indices of SIFT subcells containing \mathbf{y} . See Figure 1 for an illustration.

For technical reasons, the statistics that are used in the two proposed models are slightly different: the MaxEnt model of section 3 works on orientations at scale 0, whereas the MS-Poisson model of section 4 works on orientations computed at multiple scales. For that reason, we postpone to the next two sections the definition of the extracted statistics.

3. Stochastic models for gradient orientations. In this section, we propose a model for generating random images constrained to have prescribed local HOGs in the SIFT subcells. When designing such a model, the main difficulty arises from the fact that several SIFT subcells can overlap, and thus one has to combine the information available in all corresponding local HOGs in a way that finally complies with all the statistical constraints. In order to cope with this issue, we exploit the framework of exponential distributions to design stochastic orientation models with prescribed statistical features. The obtained distribution is “as uniform

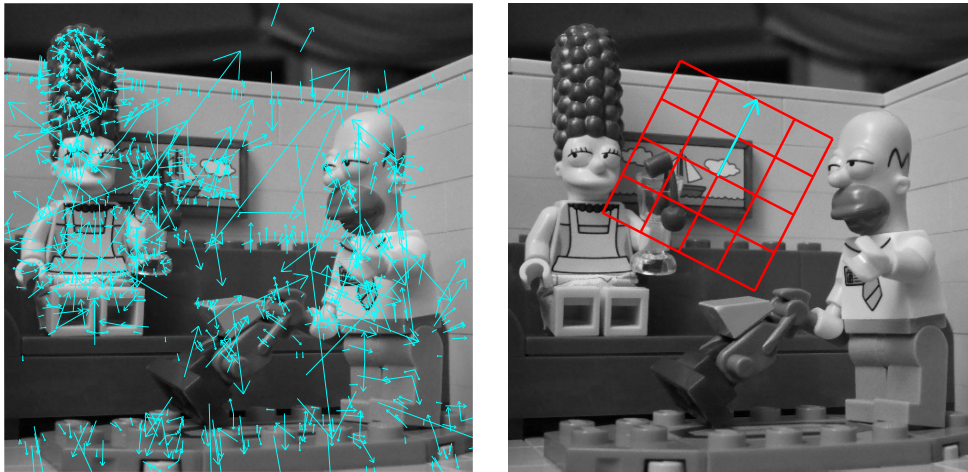


Figure 1. Examples of SIFT keypoints and subcells. *On the left, one can see an original image (courtesy of J. Delon) with superimposed SIFT oriented keypoints $(\mathbf{x}, \sigma, \alpha)$ represented by arrows originating from \mathbf{x} , with orientation α and length 6σ . On the right, we display the 16 SIFT subcells associated with one particular keypoint. Each subcell is of size $3\sigma \times 3\sigma$.*

(random) as possible” in the sense that it is of maximal entropy among all absolutely continuous distributions which satisfy the desired constraints. We combine this random orientation field with a deterministic magnitude (which is computed with the scales of locally available keypoints) in order to obtain a random objective vector field for the gradient. Finally we solve a Poisson reconstruction problem in order to get back a random image whose gradient is as close as possible as the randomly sampled objective vector field.

3.1. Exponential models with local HOG. Recall that $\Omega \subset \mathbb{Z}^2$ is a discrete rectangle. We will denote by $\mathbb{T} = \mathbb{R}/2\pi\mathbb{Z}$ the set of angles, and by \mathbb{T}^Ω the set of all possible orientation fields $\theta = (\theta(\mathbf{x}))_{\mathbf{x} \in \Omega}$ on Ω .

Extracted statistics. For simplicity, in contrast with the usual SIFT method, in this section we only extract gradient orientations at scale 0 and in addition adopt the same quantization bins for all SIFT subcells,

$$(3) \quad B_\ell = \left[(\ell - 1)\frac{\pi}{4}, \ell\frac{\pi}{4} \right) \quad (1 \leq \ell \leq 8)$$

(i.e., we do not adapt quantization to the principal orientation of the keypoint).

For all $j \in \mathcal{J}$ and $1 \leq \ell \leq 8$, we thus consider the real-valued function defined on orientation fields by

$$(4) \quad \forall \theta \in \mathbb{T}^\Omega, \quad f_{j,\ell}(\theta) = \frac{1}{|s_j|} \sum_{\mathbf{x} \in s_j} \mathbf{1}_{B_\ell}(\theta(\mathbf{x})).$$

Thus $f_{j,\ell}(\theta)$ is the proportion of points $\mathbf{x} \in s_j$ having their orientation $\theta(\mathbf{x})$ in B_ℓ .

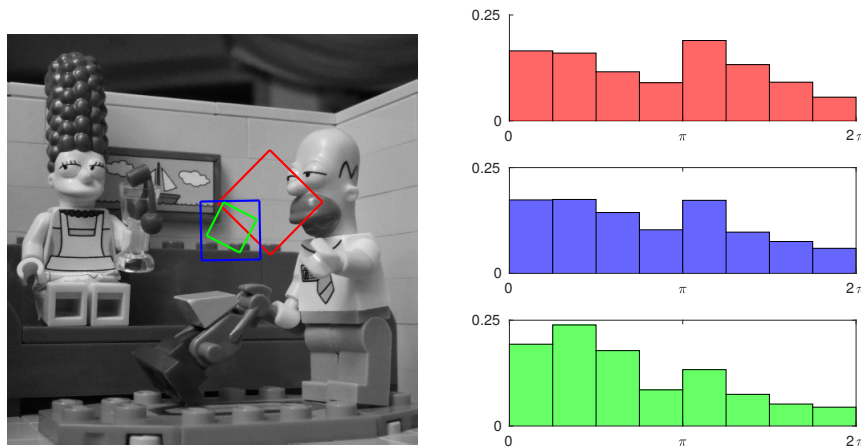


Figure 2. Extracting HOG in SIFT subcells. *On the left, we display an original image (courtesy of J. Delon) with three superimposed SIFT subcells s_j , and on the right, we display the corresponding HOG $(f_{j,\ell}(\theta_0))_{1 \leq \ell \leq 8}$ extracted in these subcells. The MaxEnt model is a probability distribution on orientation fields that will preserve on average the local HOG extracted in the SIFT subcells.*

Maximum entropy distribution. We are then interested in probability distributions P on \mathbb{T}^Ω such that

$$(5) \quad \forall j \in \mathcal{J}, \forall \ell \in \{1, \dots, 8\}, \quad \mathbb{E}_P(f_{j,\ell}(\Theta)) = \frac{1}{|s_j|} \sum_{\mathbf{x} \in s_j} \mathbb{P}(\theta(\mathbf{x}) \in B_\ell) = f_{j,\ell}(\theta_0),$$

where $\theta_0 = \text{Angle}(\nabla u_0)$ is the orientation field of the original image u_0 , and where Θ is a random orientation field with probability distribution P . In other words, we look for a random model on orientation fields which preserves on average the extracted statistics in the SIFT subcells; see Figure 2.

Let us emphasize here that we only aim at *average preservation* of the extracted statistics $(f_{j,\ell})$ because of the statistical nature of the SIFT descriptors. As will be clarified with the expression of the MaxEnt model (in particular in the case of nonoverlapping SIFT subcells), this average preservation guarantee is sufficient to precisely set the gradient orientation distribution at each point.

There are many probability distributions P on \mathbb{T}^Ω that satisfy (5), and we will be mainly interested in those that are at the same time as “random” as possible, in the sense that they are of maximal entropy. The following theorem shows the existence of such maximal entropy distributions.

Theorem 2. *There exists a family of numbers $\lambda = (\lambda_{j,\ell})_{j \in \mathcal{J}, 1 \leq \ell \leq 8}$ such that the probability distribution*

$$(6) \quad dP_\lambda = \frac{1}{Z_\lambda} \exp \left(- \sum_{j,\ell} \lambda_{j,\ell} f_{j,\ell}(\theta) \right) d\theta,$$

where the partition function Z_λ is given by $Z_\lambda = \int_{\mathbb{T}^\Omega} \exp \left(- \sum_{j,\ell} \lambda_{j,\ell} f_{j,\ell}(\theta) \right) d\theta$, satisfies

the constraints (5) and is of maximal entropy among all absolutely continuous probability distributions with respect to the Lebesgue measure $d\theta$ on \mathbb{T}^Ω satisfying the constraints (5).

Proof. This result directly follows from the general theorem given in [46]. The only difficulty is in handling the hypothesis of linear independence of the $f_{j,\ell}$. In our framework, the $f_{j,\ell}$ are not independent (in particular because $\sum_{\ell=1}^8 f_{j,\ell} = 1$, and also because there may be other dependencies—for instance, when one subcell is exactly the union of two smaller subcells). But one can still apply the theorem to an extracted linearly independent subfamily. This gives the existence of the solution for the initial family $(f_{j,\ell})$ (but of course not the unicity). ■

Remark. We do not repeat here the argument (based on Lagrange multipliers) showing that maximizing entropy under constraints (5) leads to exponential distributions. However, once a solution P_λ has been computed, and if P is an absolutely continuous probability distribution satisfying (5), one can write the Kullback–Leibler divergence using the entropy $H(P)$:

$$(7) \quad D(P||P_\lambda) = \int \log \left(\frac{P(\theta)}{P_\lambda(\theta)} \right) P(\theta) d\theta = -H(P) + \log Z_\lambda + \sum \lambda_{j,\ell} f_{j,\ell}(\theta_0),$$

which shows that maximizing $H(P)$ under (5) is equivalent to minimize $D(P||P_\lambda)$. In particular, this shows that the maximal entropy distribution under (5) is unique (because of the strict concavity of the entropy) even if there may be several sets of parameters λ corresponding to that solution.

Independence property of the MaxEnt model.

Proposition 3. *Under P_λ the values $\Theta(\mathbf{x})$ are independent. In addition, the probability density function of $\Theta(\mathbf{x})$ is given by*

$$(8) \quad \frac{1}{Z_{\lambda,\mathbf{x}}} e^{-\varphi_{\lambda,\mathbf{x}}} = \frac{1}{Z_{\lambda,\mathbf{x}}} \sum_{\ell=1}^8 \exp \left(- \sum_{j \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{j,\ell}}{|s_j|} \right) \mathbf{1}_{B_\ell},$$

$$(9) \quad \text{where } Z_{\lambda,\mathbf{x}} = \sum_{\ell=1}^8 \exp \left(- \sum_{j \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{j,\ell}}{|s_j|} \right) |B_\ell|.$$

Proof. Taking the logarithm of (6), one can group the terms corresponding to the same pixel \mathbf{x} so that

$$(10) \quad -\log \frac{dP_\lambda}{d\theta} - \log Z_\lambda = \sum_{j \in \mathcal{J}, 1 \leq \ell \leq 8} \lambda_{j,\ell} f_{j,\ell}(\theta) = \sum_{\mathbf{x} \in \Omega} \varphi_{\lambda,\mathbf{x}}(\theta(\mathbf{x})),$$

$$(11) \quad \text{where } \varphi_{\lambda,\mathbf{x}} = \sum_{\ell=1}^8 \left(\sum_{j \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{j,\ell}}{|s_j|} \right) \mathbf{1}_{B_\ell}.$$

We thus obtain that P_λ can be written in a separable form. ■

On the one hand, this proposition shows that for a given λ , one can easily sample from the model P_λ . On the other hand, it also allows us to compute several statistics associated with this model. In particular, we can compute for any bounded measurable function $\psi : \mathbb{T} \rightarrow \mathbb{C}$

$$(12) \quad \mathbb{E}_{P_\lambda}[\psi(\Theta(\mathbf{x}))] = \frac{\sum_{\ell=1}^8 \exp\left(-\sum_{j \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{j,\ell}}{|s_j|}\right) \int_{B_\ell} \psi(t) dt}{\sum_{\ell=1}^8 \exp\left(-\sum_{j \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{j,\ell}}{|s_j|}\right) |B_\ell|}.$$

It also allows us to compute the expected value of the statistics $f(\Theta)$ in the model P_λ (which will be useful in section 3.3):

$$(13) \quad \mathbb{E}_{P_\lambda}[f_{j,\ell}(\Theta)] = \frac{1}{|s_j|} \sum_{\mathbf{x} \in s_j} \mathbb{P}(\Theta(\mathbf{x}) \in B_\ell) = \frac{1}{|s_j|} \sum_{\mathbf{x} \in s_j} \frac{\exp\left(-\sum_{k \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{k,\ell}}{|s_k|}\right) |B_\ell|}{\sum_{1 \leq \ell' \leq 8} \exp\left(-\sum_{k \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{k,\ell'}}{|s_k|}\right) |B_{\ell'}|}.$$

But it remains to show how to estimate λ in order to satisfy the constraints (5). These constraints can be rewritten as

$$(14) \quad \forall j, \ell, \quad \sum_{\mathbf{x} \in s_j} \frac{1}{Z_{\lambda,\mathbf{x}}} \exp\left(-\sum_{k \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{k,\ell}}{|s_k|}\right) |B_\ell| = |\{\mathbf{x} \in s_j ; \theta_0(\mathbf{x}) \in B_\ell\}|.$$

Notice that this system is highly nonlinear and is in general difficult to solve.

A simple case: Nonoverlapped SIFT subcells. When a SIFT subcell s_j is not overlapped, we have, for any $\mathbf{x} \in s_j$, $|\mathcal{J}(\mathbf{x})| = 1$ and therefore

$$(15) \quad Z_{\lambda,\mathbf{x}} = \sum_{\ell=1}^8 \exp\left(-\frac{\lambda_{j,\ell}}{|s_j|}\right) |B_\ell|.$$

Then (14) gives

$$(16) \quad \forall \ell, \quad \frac{1}{Z_{\lambda,\mathbf{x}}} \exp\left(-\frac{\lambda_{j,\ell}}{|s_j|}\right) = \frac{|\{\mathbf{x} \in s_j ; \theta_0(\mathbf{x}) \in B_\ell\}|}{|s_j| |B_\ell|} = f_{j,\ell}(\theta_0),$$

which gives the marginal distribution on any $\mathbf{x} \in s_j$:

$$(17) \quad \frac{1}{Z_{\lambda,\mathbf{x}}} e^{-\varphi_{\lambda,\mathbf{x}}} = \sum_{\ell=1}^8 \frac{|\{\mathbf{x} \in s_j ; \theta_0(\mathbf{x}) \in B_\ell\}|}{|s_j| |B_\ell|} \mathbf{1}_{B_\ell} = \sum_{\ell=1}^8 f_{j,\ell}(\theta_0) \frac{1}{|B_\ell|} \mathbf{1}_{B_\ell}.$$

So when the subcells do not overlap, the maximum entropy distribution only amounts to independent resampling of the local HOGs, as expected. Notice that we indeed obtain a unique maximal entropy distribution. However, the solutions λ are only unique up to the addition of a constant: indeed the last calculation shows that for a nonoverlapped subcell s_j , there exists a constant $c_j > 0$ such that

$$(18) \quad \forall \ell, \quad \lambda_{j,\ell} = -|s_j|(\log f_{j,\ell}(\theta_0) + \log c_j).$$

Maximum-likelihood estimation. If the SIFT subcells intersect, there is no longer an explicit solution. To cope with that, as in [69] we use a numerical scheme to find the maximum entropy distribution P_λ . The solution can be obtained with a traditional maximum likelihood estimation technique, as will be detailed here. Indeed, the minus-log-likelihood function can be written as

$$(19) \quad \Phi(\lambda) = \log Z_\lambda + \sum_{j,\ell} \lambda_{j,\ell} f_{j,\ell}(\theta_0).$$

The gradient of Φ can be obtained by differentiating the partition function

$$(20) \quad \frac{\partial \log Z_\lambda}{\partial \lambda_{j,\ell}} = \frac{1}{Z_\lambda} \frac{\partial Z_\lambda}{\partial \lambda_{j,\ell}} = -\mathbb{E}_{P_\lambda} [f_{j,\ell}(\Theta)],$$

which gives

$$(21) \quad \frac{\partial \Phi}{\partial \lambda_{j,\ell}} = f_{j,\ell}(\theta_0) - \mathbb{E}_{P_\lambda} [f_{j,\ell}(\Theta)].$$

Notice that $\nabla \Phi(\lambda) = 0$ if and only if P_λ satisfies the constraints (5).

Similarly, we can also obtain the second order derivatives

$$(22) \quad \frac{\partial^2 \Phi}{\partial \lambda_{j,\ell} \partial \lambda_{j',\ell'}} = \mathbb{E}_{P_\lambda} \left[(f_{j,\ell}(\Theta) - \mathbb{E}_{P_\lambda} [f_{j,\ell}(\Theta)]) (f_{j',\ell'}(\Theta) - \mathbb{E}_{P_\lambda} [f_{j',\ell'}(\Theta)]) \right].$$

One can observe that this Hessian matrix $\nabla^2 \Phi(\lambda)$ is actually the covariance of the vector $f(\Theta)$ when Θ has distribution P_λ . In particular it is a semipositive definite matrix, which shows that Φ is a convex function. The global minima of Φ are exactly the points λ where $\nabla \Phi$ vanishes, which is equivalent to having the constraints (5) on P_λ .

Therefore, we can compute the solution P_λ by a gradient descent algorithm in order to minimize Φ . The complete algorithm is summarized in section 3.3. Since Φ is not strictly convex, we will not have a guarantee of convergence on the iterates, but on the function values. Since $|f_{j,\ell}(\theta)| \leq 1$, it is straightforward to see that all coefficients of the Hessian $\nabla^2 \Phi(\lambda)$ have modulus ≤ 1 . Therefore, the ℓ^2 operator norm of $\nabla^2 \Phi$ is bounded by $8|\mathcal{J}|$, which implies that $\nabla \Phi$ is L -Lipschitz with $L = 8|\mathcal{J}|$. Denoting by λ^k the iterates of the gradient descent with constant step size $h < \frac{2}{L}$, [48, Thm. 2.1.14] gives

$$(23) \quad \Phi(\lambda^k) - \min \Phi = \mathcal{O}\left(\frac{1}{k}\right).$$

Let us also mention that since Φ is convex smooth, it would be possible to use higher order optimization schemes to minimize Φ . However, Newton's method will be in general too costly because of the dimension of the system and because the Hessian may be ill-conditioned.

3.2. Monoscale Poisson reconstruction. Now that we have built a random orientation field Θ with maximum entropy distribution P_λ , we will use it to propose a target vector field V for the image gradient. More precisely, we set the gradient magnitude at \mathbf{x} in a deterministic manner, as the inverse scale of the smallest subcell that covers \mathbf{x} . For pixels \mathbf{x} which lie

outside the SIFT subcells, we set $V(\mathbf{x}) = 0$. This choice allows us to give more weight to the locations for which we have information at a finer scale. It is also motivated by the following homogeneity argument. Assume that $u : \mathbb{R}^2 \rightarrow \mathbb{R}$ has a keypoint (\mathbf{x}, σ) , and for $a > 0$ let $v(\mathbf{y}) = u(\frac{\mathbf{y}}{a})$. Then, thanks to Proposition 1, v has a keypoint $(a\mathbf{x}, a\sigma)$. Let us compare the mean gradient magnitude at scale σ in the corresponding subcell s to the analogous quantity for v . A simple computation shows that

$$\frac{1}{|as|} \int_{\lambda_s} |\nabla g_{a\sigma} * v(\mathbf{y})| d\mathbf{y} = \frac{1}{a} \frac{1}{|s|} \int_s |\nabla g_\sigma * u(\mathbf{y})| d\mathbf{y},$$

so that the mean gradient magnitude in the subcell is multiplied by $\frac{1}{a}$ with the change of scale. From this calculation we get the following remark: if two very similar shapes (with similar gray levels) are seen in the image at two different scales with ratio a , then we can obtain a pairwise matching of their SIFT keypoints, and the ratio between the mean gradient magnitude of the two matched subcells is $1/a$. Of course this remark does not extend to the comparison of two SIFT subcells with very different geometric content, but it still provides a general rule for fixing the gradient magnitude as the inverse of the scale. Therefore, we get the random objective vector field

$$(24) \quad \forall \mathbf{x} \in \Omega, \quad V(\mathbf{x}) = \left(\max_{j \in \mathcal{J}(\mathbf{x})} \frac{1}{\sigma_j} \right) e^{i\Theta(\mathbf{x})} \mathbf{1}_{\mathcal{J}(\mathbf{x}) \neq \emptyset}.$$

The aim of the Poisson reconstruction is to compute an image whose gradient is as close as possible to the vector field $V = (V_1, V_2)$. In the case of image editing, this technique has been proposed by Pérez, Gangnet, and Blake [51] in order to copy pieces of an image into another one in a seamless way. More precisely, the goal is to minimize the functional

$$(25) \quad F(u) = \sum_{\mathbf{x} \in \Omega} \|\nabla u(\mathbf{x}) - V(\mathbf{x})\|_2^2.$$

Since $F(c + u) = F(u)$ for any constant c , we can impose $\sum_{\mathbf{x} \in \Omega} u(\mathbf{x}) = 0$. Thus we set

$$(26) \quad U = \text{Argmin} \left\{ F(u); u : \Omega \rightarrow \mathbb{R} \text{ and such that } \sum_{\mathbf{x} \in \Omega} u(\mathbf{x}) = 0 \right\}.$$

If we use periodic boundary conditions for the gradient, we can solve this problem with the discrete Fourier transform [43]. Indeed, if we use the simple derivation scheme based on periodic convolutions

$$(27) \quad \nabla u(\mathbf{x}) = \begin{pmatrix} \partial_1 * u(\mathbf{x}) \\ \partial_2 * u(\mathbf{x}) \end{pmatrix}, \quad \text{where} \quad \begin{cases} \partial_1 &= \delta_{(0,0)} - \delta_{(1,0)}, \\ \partial_2 &= \delta_{(0,0)} - \delta_{(0,1)}, \end{cases}$$

the problem can be expressed in the Fourier domain with Parseval formula since

$$(28) \quad F(u) = \frac{1}{|\Omega|} \sum_{\boldsymbol{\xi}} |\widehat{\partial_1}(\boldsymbol{\xi})\widehat{u}(\boldsymbol{\xi}) - \widehat{V}_1(\boldsymbol{\xi})|_2^2 + |\widehat{\partial_2}(\boldsymbol{\xi})\widehat{u}(\boldsymbol{\xi}) - \widehat{V}_2(\boldsymbol{\xi})|_2^2.$$

Thus, for each ξ we have a barycenter problem which is simply solved by

$$(29) \quad \forall \xi \neq 0, \quad \widehat{U}(\xi) = \frac{\widehat{\partial}_1(\xi)\widehat{V}_1(\xi) + \widehat{\partial}_2(\xi)\widehat{V}_2(\xi)}{|\widehat{\partial}_1(\xi)|^2 + |\widehat{\partial}_2(\xi)|^2} \quad \text{and} \quad \widehat{U}(0) = 0.$$

Let us emphasize (with the capital letter U) that the solution of this problem is random because the target field V is random.

Using the notation $\nabla = (\partial_1, \partial_2)^T$, $\widehat{\nabla} = (\widehat{\partial}_1, \widehat{\partial}_2)^T$, $z^* = \bar{z}^T$, we can write

$$(30) \quad \widehat{U}(\xi) = \widehat{\nu}(\xi)\widehat{V}(\xi), \quad \text{where} \quad \widehat{\nu}(\xi) = \begin{cases} \frac{\widehat{\nabla}(\xi)^*}{|\widehat{\nabla}(\xi)|^2} & \text{if } \xi \neq 0, \\ 0 & \text{if } \xi = 0. \end{cases}$$

Notice that $\widehat{\nu}(\xi) \in \mathbb{C}^{1 \times 2}$ and $\widehat{V}(\xi) \in \mathbb{C}^{2 \times 1}$, so that (30) is equivalent to

$$(31) \quad U = \nu * V = \nu_1 * V_1 + \nu_2 * V_2.$$

In other words, ν is the (vector-valued) convolution kernel associated to the Poisson reconstruction. This expression allows us to compute the moments of the random field U (see also section 4.3 for a detailed look at a more general calculation).

3.3. Algorithm. Here we summarize the algorithm for estimating and sampling the MaxEnt model proposed in this section. In Figure 3 we display an example of reconstruction with the MaxEnt model.

Algorithm: Estimating and Sampling the MaxEnt Model

- Maximum-likelihood estimation of λ
 - Compute the observed statistics $f(\theta_0) = (f_{j,\ell}(\theta_0))_{j,\ell}$.
 - Initialization $\lambda \leftarrow 0$. Choose a step size $h < \frac{4}{|\mathcal{J}|}$.
 - For $N(= 10000)$ iterations, compute $\bar{f} = \mathbb{E}_{P_\lambda}[f]$ using (13) and set

$$\lambda \leftarrow \lambda - h(f(\theta_0) - \bar{f}).$$

- Draw a sample θ according to the distribution P_λ .
- Compute the corresponding target vector field

$$(32) \quad V(\mathbf{x}) = \left(\max_{j \in \mathcal{J}(\mathbf{x})} \frac{1}{\sigma_j} \right) e^{i\theta(\mathbf{x})} \mathbf{1}_{\mathcal{J}(\mathbf{x}) \neq \emptyset}.$$

- Compute a sample u of MaxEnt via the Poisson reconstruction (29).
-

For images having many SIFT keypoints in overlapping positions, this algorithm may be slow to converge, as can be observed on the case of Figure 3. This case is relatively simple because it has only 187 keypoints, but this corresponds already to $8 \times 16 \times 187 \approx 24000$ $\lambda_{j,\ell}$ parameters to estimate. This is why we use a stopping criterion based on a maximal number of iterations.

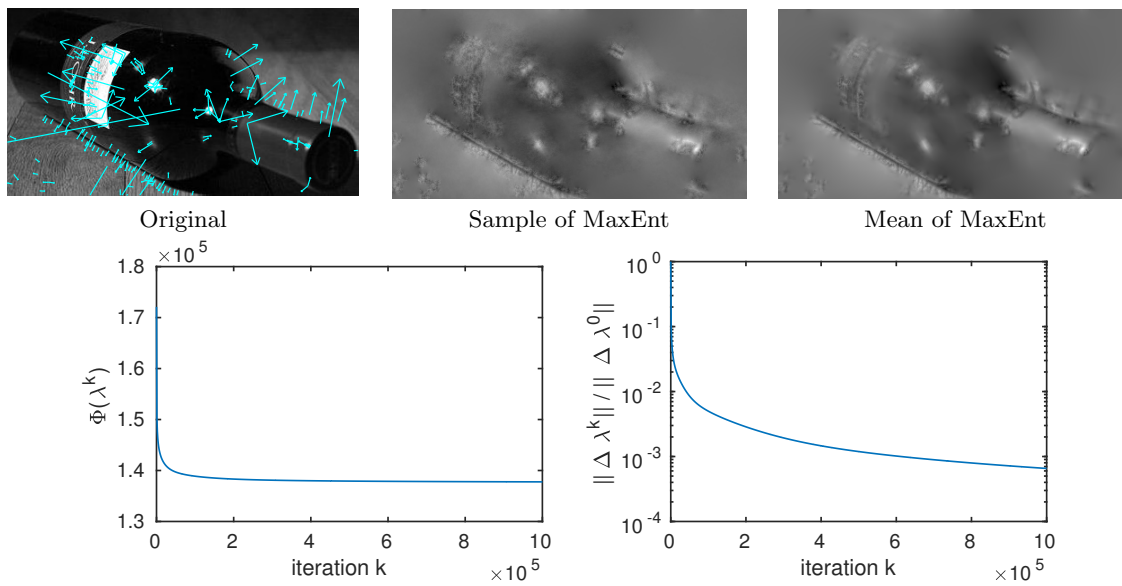


Figure 3. Reconstruction with the MaxEnt model. In the first row from left to right, we display an original image with 187 oriented keypoints superimposed, a sample of the associated MaxEnt model, and the expectation of the MaxEnt model. In the second row we display the evolution of Φ along the iterates, and also the behavior of the difference between iterates $\Delta\lambda^k = \lambda^k - \lambda^{k-1}$. The value of Φ stabilizes in about 10^5 iterations. One can remark that both reconstructions show several important structures of the original image. The mean reconstruction is of course smoother than a sample of the model (because pixels are sampled independently; see Proposition 3).

3.4. Discussion on MaxEnt model. One drawback of MaxEnt is that the guarantee on the local distributions of orientations is lost after the Poisson reconstruction step. One way to cope with that would be to consider a model that operates directly on the image values, and not on the orientation field. Theorem 2 could be extended to statistics like

$$(33) \quad \tilde{f}_{j,\ell}(u) = \frac{1}{|s_j|} \sum_{\mathbf{x} \in s_j} \mathbf{1}_{B_\ell}(\text{Angle}(\nabla u(\mathbf{x}))).$$

It is even possible to consider multiscale statistics using $\nabla g_{\sigma_j} * u$ instead of ∇u (as will be the case in section 4). But the analogue of Proposition 3 would not hold for these models, so that sampling should rely on a Gibbs strategy. Its cost would be clearly prohibitive in the multiscale case due to the large Markov neighborhood size. Even in the monoscale case the convergence of this Gibbs sampler may be very long depending on the parameters λ ; and since we would need one sample per iteration of gradient descent to estimate λ , we chose to leave it aside and concentrate on models with reasonably fast sampling.

Also, one can consider another orientation model in which the local HOGs are computed with a quantization that depends on the keypoint orientation. The independence property still holds for this model, and the marginal orientations still have a piecewise constant density, but the number of parameters would be much larger (there would be as many ℓ 's as bins of a subdivision that is adapted to all keypoints orientations). Therefore this model is practically intractable, and also only of minor interest. Indeed, in view of the results of Figure 3, it is

likely that the used quantization has only a minor impact on the visual results (provided that we still have a minimal number of bins).

4. Multiscale Poisson model. In this section, we propose a stochastic model, called MS-Poisson, for reconstruction using multiscale local HOGs computed in SIFT subcells. This new model is based on a heuristic algorithm for orientation resampling in all SIFT subcells. Therefore, in contrast to the MaxEnt model, the MS-Poisson model can be straightforwardly sampled using the multiscale local HOGs and does not require an iterative estimation procedure. Another difference is that MS-Poisson is designed to combine information at multiple scales, whereas MaxEnt only operates with the gradient at scale 0.

4.1. Construction of MS-Poisson model.

Extracted statistics. The MS-Poisson model is based on local statistics on multiscale gradient orientations. More precisely, in s_j we extract the quantized HOG at scale σ_j :

$$(34) \quad H_{j,\ell} = \frac{1}{|s_j|} |\{\mathbf{x} \in s_j ; \text{Angle}(\nabla g_{\sigma_j} * u_0)(\mathbf{x}) - \alpha_j \in [(\ell - 1)\frac{\pi}{4}, \ell\frac{\pi}{4}]\}|.$$

In view of resampling, this local HOG can be identified to a piecewise constant density function:

$$(35) \quad h_j = \frac{4}{\pi} \sum_{\ell=1}^8 H_{j,\ell} \mathbf{1}_{[\alpha_j + (\ell-1)\frac{\pi}{4}, \alpha_j + \ell\frac{\pi}{4}]}.$$

Notice that, in contrast to the statistics (4) used in the MaxEnt model, the quantization here depends on the local orientation α_j .

Target vector fields at multiple scales. Using the local orientation distributions h_j , we define vector fields $V_j : \Omega \rightarrow \mathbb{R}^2$ that will serve as objective gradients at scale σ_j in the SIFT subcell s_j . We propose setting

$$(36) \quad \forall \mathbf{x} \in \Omega, \quad V_j(\mathbf{x}) = \frac{1}{\sigma_j} e^{i\gamma_j(\mathbf{x})} \mathbf{1}_{s_j}(\mathbf{x}),$$

where the orientations $\gamma_j(\mathbf{x})$ are independently sampled according to the distribution h_j . Again, as justified in section 3.2, we set the gradient magnitude in a deterministic way using the inverse of the scale σ_j . Once these vector fields V_j have been sampled, we obtain an image U by solving a multiscale Poisson problem, as explained in the next subsection.

4.2. Multiscale Poisson reconstruction. In order to simultaneously constrain the gradient at several scales $(\sigma_j)_{j \in \mathcal{J}}$, we propose considering the following multiscale Poisson energy:

$$(37) \quad G(u) = \sum_{j \in \mathcal{J}} w(\sigma_j) \sum_{\mathbf{x} \in \Omega} \|\nabla(g_{\sigma_j} * u)(\mathbf{x}) - V_j(\mathbf{x})\|_2^2,$$

where g_σ is the Gaussian kernel of standard deviation σ , $V_j = (V_{j,1}, V_{j,2})^T$ is the objective gradient at scale σ_j , and $\{w(\sigma_j), j \in \mathcal{J}\}$ is a set of weights. In our application, since there are

more keypoints in the fine scales (i.e., with small σ_j), and since the keypoints at fine scales are generally more informative, a reasonable choice is to take all weights $w(\sigma_j) = 1$. But we keep these weights in the formula for the sake of generality. We thus set

$$(38) \quad U = \operatorname{Argmin} \left\{ G(u); u : \Omega \rightarrow \mathbb{R} \text{ and such that } \sum_{\mathbf{x} \in \Omega} u(\mathbf{x}) = 0 \right\}.$$

Again, with periodic boundary conditions, this problem can be expressed in the Fourier domain as

$$(39) \quad G(u) = \frac{1}{|\Omega|} \sum_{j \in \mathcal{J}} \sum_{\boldsymbol{\xi}} w(\sigma_j) \left(|\widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \widehat{\partial}_1(\boldsymbol{\xi}) \widehat{u}(\boldsymbol{\xi}) - \widehat{V}_{j,1}(\boldsymbol{\xi})|_2^2 + |\widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \widehat{\partial}_2(\boldsymbol{\xi}) \widehat{u}(\boldsymbol{\xi}) - \widehat{V}_{j,2}(\boldsymbol{\xi})|_2^2 \right).$$

As for the monoscale Poisson problem, the solution U is still a barycenter given by $\widehat{U}(0) = 0$ and

$$(40) \quad \forall \boldsymbol{\xi} \neq 0, \quad \widehat{U}(\boldsymbol{\xi}) = \frac{\sum_{j \in \mathcal{J}} w(\sigma_j) \widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \left(\overline{\widehat{\partial}_1(\boldsymbol{\xi})} \widehat{V}_{j,1}(\boldsymbol{\xi}) + \overline{\widehat{\partial}_2(\boldsymbol{\xi})} \widehat{V}_{j,2}(\boldsymbol{\xi}) \right)}{\sum_{j \in \mathcal{J}} w(\sigma_j) |\widehat{g}_{\sigma_j}(\boldsymbol{\xi})|^2 \left(|\widehat{\partial}_1(\boldsymbol{\xi})|^2 + |\widehat{\partial}_2(\boldsymbol{\xi})|^2 \right)}.$$

Let us remark that in the above formula, we have $\widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \in \mathbb{R}$ since g_{σ_j} is even.

Regularization. Notice that, depending on the finest scale, the denominator may numerically vanish in the high frequencies because of the term $\widehat{g}_{\sigma_j}(\boldsymbol{\xi})$ (as it is the case in a deconvolution problem). Therefore, it may be useful to add a regularization term controlled by a parameter $\mu > 0$. Then, if we set

$$(41) \quad U = \operatorname{Argmin} \left\{ G(u) + \mu \|\nabla u\|_2^2; u : \Omega \rightarrow \mathbb{R} \text{ and such that } \sum_{\mathbf{x} \in \Omega} u(\mathbf{x}) = 0 \right\},$$

then we get the well-defined solution U given by $\widehat{U}(0) = 0$ and

$$(42) \quad \forall \boldsymbol{\xi} \neq 0, \quad \widehat{U}(\boldsymbol{\xi}) = \frac{\sum_{j \in \mathcal{J}} w(\sigma_j) \widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \left(\overline{\widehat{\partial}_1(\boldsymbol{\xi})} \widehat{V}_{j,1}(\boldsymbol{\xi}) + \overline{\widehat{\partial}_2(\boldsymbol{\xi})} \widehat{V}_{j,2}(\boldsymbol{\xi}) \right)}{\left(\mu + \sum_{j \in \mathcal{J}} w(\sigma_j) |\widehat{g}_{\sigma_j}(\boldsymbol{\xi})|^2 \right) \left(|\widehat{\partial}_1(\boldsymbol{\xi})|^2 + |\widehat{\partial}_2(\boldsymbol{\xi})|^2 \right)}.$$

As we will see in section 5.1, the parameter μ allows us to attenuate the noise generated by the randomly sampled gradient fields in the fine-scale SIFT subcells. We will see (empirically) that the value $\mu = 50$ realizes a good compromise between recovered details and smoothness.

We end this subsection by summarizing the MS-Poisson sampling algorithm.

Algorithm: Sampling the MS-Poisson Model

- In each subcell s_j , draw independent orientations $\gamma_j(\mathbf{x}), \mathbf{x} \in s_j$ according to the p.d.f. h_j .
 - Set $V_j = \frac{1}{\sigma_j} \mathbf{1}_{s_j} e^{i\gamma_j}$.
 - Compute U by solving the MS-Poisson problem (41) with targets V_j , with $w(\sigma_j) = 1$ and $\mu = 50$.
-

Remark. In (42), one can observe that the solution to MS-Poisson actually solves a mono-scale Poisson problem with objective vector field V whose Fourier transform is given by

$$(43) \quad \widehat{V}(\boldsymbol{\xi}) = \frac{\sum_{j \in \mathcal{J}} w(\sigma_j) \widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \widehat{V}_j(\boldsymbol{\xi})}{\mu + \sum_{j \in \mathcal{J}} w(\sigma_j) |\widehat{g}_{\sigma_j}(\boldsymbol{\xi})|^2}.$$

4.3. First and second order moments. In order to compute the statistics of the MS-Poisson model, we remark that the multiscale Poisson reconstruction is actually a linear process. Indeed, for each j , let $\nu_j : \Omega \rightarrow \mathbb{R}^{1 \times 2}$ be the vector-valued kernel defined by its discrete Fourier transform

$$(44) \quad \forall \boldsymbol{\xi} \neq 0, \quad \widehat{\nu}_j(\boldsymbol{\xi}) = \frac{w(\sigma_j) \widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \widehat{V}(\boldsymbol{\xi})^*}{\left(\mu + \sum_{j' \in \mathcal{J}} w(\sigma_{j'}) |\widehat{g}_{\sigma_{j'}}(\boldsymbol{\xi})|^2 \right) |\widehat{V}(\boldsymbol{\xi})|^2} \quad \text{and } \widehat{\nu}_j(0) = 0.$$

Then, as in section 3.2, we get the convolutive expression

$$(45) \quad U = \sum_{j \in \mathcal{J}} \nu_j * V_j = \sum_{j \in \mathcal{J}} \left(\nu_{j,1} * V_{j,1} + \nu_{j,2} * V_{j,2} \right).$$

From this expression we can compute the moments of U . By linearity

$$(46) \quad \mathbb{E}(U) = \sum_{j \in \mathcal{J}} \nu_j * \mathbb{E}(V_j),$$

so that computing this expectation only amounts to computing $\mathbb{E}(V_j) = \frac{1}{\sigma_j} \mathbf{1}_{s_j} \mathbb{E}(e^{i\gamma_j})$.

We can also compute the variance. Since the objective fields $(V_j)_{j \in \mathcal{J}}$ are independent, we have

$$(47) \quad \text{Var}(U(\mathbf{x})) = \sum_{j \in \mathcal{J}} \text{Var}(\nu_j * V_j(\mathbf{x})).$$

Also, the $V_j(\mathbf{y})$ for different pixels \mathbf{y} are independent, so that

$$(48) \quad \text{Var}(\nu_j * V_j(\mathbf{x})) = \text{Var}\left(\sum_{\mathbf{y} \in \Omega} \nu_j(\mathbf{x} - \mathbf{y})V_j(\mathbf{y})\right) = \sum_{\mathbf{y} \in \Omega} \text{Var}(\nu_j(\mathbf{x} - \mathbf{y})V_j(\mathbf{y}))$$

$$(49) \quad = \sum_{\mathbf{y} \in \Omega} \nu_j(\mathbf{x} - \mathbf{y})\text{Cov}(V_j(\mathbf{y}))\nu_j^T(\mathbf{x} - \mathbf{y})$$

$$(50) \quad = \sum_{\mathbf{y} \in \Omega} \nu_{j,1}^2(\mathbf{x} - \mathbf{y})\text{Var}(V_{j,1}(\mathbf{y})) + \nu_{j,2}^2(\mathbf{x} - \mathbf{y})\text{Var}(V_{j,2}(\mathbf{y}))$$

$$(51) \quad + 2\nu_{j,1}(\mathbf{x} - \mathbf{y})\nu_{j,2}(\mathbf{x} - \mathbf{y})\text{Cov}(V_{j,1}(\mathbf{y}), V_{j,2}(\mathbf{y})).$$

Therefore the variance of this model can be obtained by summing convolutions of the kernels ν_j with the covariances of V_j . Since $V_j(\mathbf{y}) = \frac{1}{\sigma_j} e^{i\gamma_j(\mathbf{y})} \mathbf{1}_{s_j}$, where $\gamma_j(\mathbf{y})$ has p.d.f. h_j given by (34), we can explicitly compute its covariance.

More generally, we can compute the covariance between two pixel values of U in a similar way, which gives

$$(52) \quad \text{Cov}(U(\mathbf{x}), U(\mathbf{y})) = \sum_{j \in \mathcal{J}} \sum_{\mathbf{z} \in \Omega} \nu_j(\mathbf{x} - \mathbf{z})\text{Cov}(V_j(\mathbf{z}))\nu_j^T(\mathbf{y} - \mathbf{z}).$$

5. Results and discussion. In this section, we give empirical evidence that both models MS-Poisson and MaxEnt are able to generate images that are similar to the original image in many aspects, which is confirmed by several quantitative results (in particular based on normalized correlations). We discuss the impact of the regularization parameter μ of the MS-Poisson model on the quality of the sampled images. We also compare MaxEnt and MS-Poisson in terms of local variance of the sampled images, and also in terms of resulting SIFT keypoints computed in the sampled images. After explaining how to adapt the MS-Poisson model to operate on true SIFT descriptors we compare with previous approaches of [64, 19]. Finally we discuss the impact of the keypoints definition on the quality of the reconstruction.

5.1. Results with MaxEnt and MS-Poisson model.

5.1.1. Comparison between MaxEnt and MS-Poisson. Let us first compare the reconstruction results obtained with MaxEnt and with MS-Poisson. In Figure 4, using an original image with 386 keypoints, we display a sample of MaxEnt and a sample of MS-Poisson, together with the expected images of these models. One first remark is that both models are able to retrieve several geometric structures of the original image, so that much of the semantic content of the image can still be understood. For both models, one can observe that the samples are very close to the expected image, which will be later confirmed by the variance analysis in Figure 8.

One crucial difference between MaxEnt and MS-Poisson is that they do not rely on the same gradient information. Indeed, MS-Poisson exploits gradients extracted at multiple scales, while MaxEnt only operates with gradients at scale $\sigma = 0$ (i.e., the same scale as the image). This is why the results obtained with MS-Poisson will generally look blurrier than those obtained with MaxEnt. Also, because of the multiscale nature of the input of MS-Poisson, the corresponding optimization problem had to be regularized, and the adopted H^1 -regularization

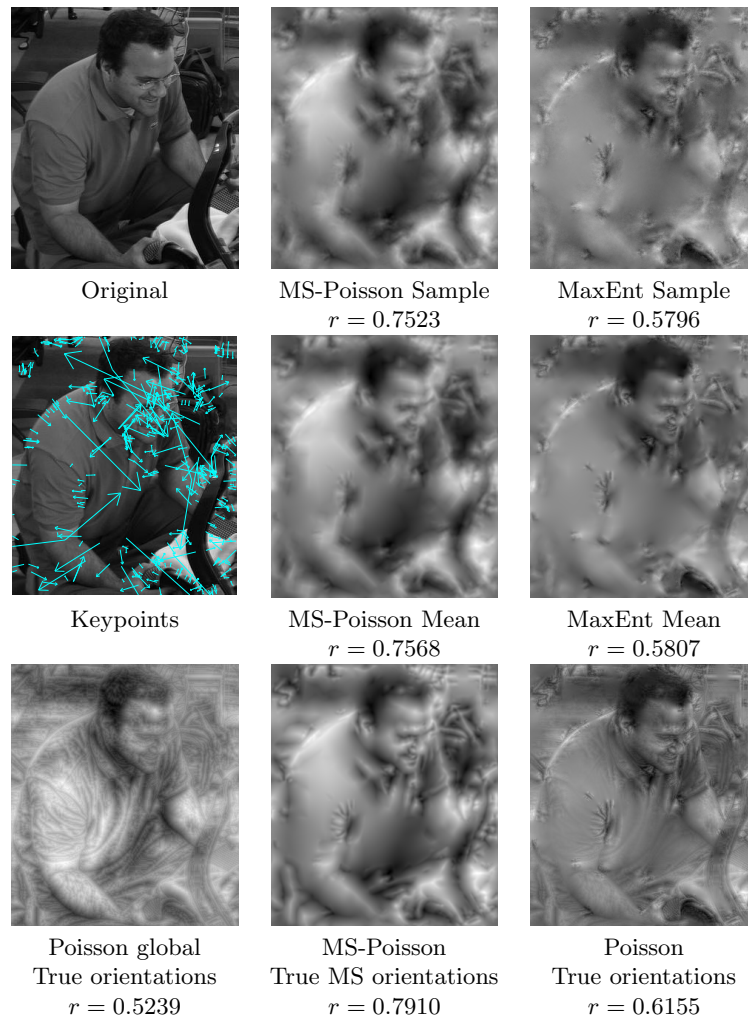


Figure 4. Reconstruction results with MaxEnt and MS-Poisson models. In the first column we display an original image, the corresponding oriented keypoints, and the Poisson reconstruction with true gradient orientations of the whole image and magnitude set to 1. In the second column we display a sample of the MS-Poisson model, the expectation of this model, and the multiscale Poisson reconstruction using the true multiscale gradient orientations in the SIFT subcells. In the third column, we display a sample of the MaxEnt model, the expectation of this model, and the Poisson reconstruction using the true gradient orientations in the SIFT subcells. For each result we indicate the value of the normalized correlation r with respect to the original image. See the text for comments on these results. (Images are better seen in the electronic version.)

term is a source of blur in the result as well. This is confirmed by Figure 5, where we display several MS-Poisson reconstructions with varying regularization parameter μ . In Figure 5 and in many other experiments, we observed that the parameter $\mu = 50$ realizes a good compromise between preserving geometric structures and removing spurious oscillations.

In the last row of Figure 4, we also compare with the reconstructions obtained with the true gradient orientations (resp., multiscale gradient orientations) computed in the SIFT subcells and the gradient magnitude computed as in MaxEnt (resp., MS-Poisson). So the difference



Figure 5. Influence of the regularization parameter μ in MS-Poisson. *As expected, increasing μ penalizes more the L^2 norm of the gradient and thus makes the image blurrier. Here again we indicate the value of the normalized correlation r with respect to the original image. We empirically observed that a good compromise between recovered details and smoothness is often attained around $\mu = 50$. (Images are better seen in the electronic version.)*

with MaxEnt (or MS-Poisson) is that local (multiscale) gradient orientations are not pooled in histograms but directly extracted pixelwise; in other words, there is no local resampling of the orientations. Thus, in some sense, these images are the best ones we could hope for using Poisson reconstruction. Comparing these images with samples of MS-Poisson and MaxEnt precisely shows the effect of local resampling of the (multiscale) orientations; observe in particular the man's face and also the folds of his shirt. These images thus correspond to much more precise reconstructions, but it is interesting to notice that in certain regions where attention will be focused (near the face, for example), there are enough keypoints at fine scales in order to get back satisfying pieces of images even after local resampling. Also, one must keep in mind that the loss of the gradient magnitude information is in practice difficult to cope with and may force us to erroneously amplify the noise in the reconstruction. As one can see in the bottom left of Figure 4, it is obvious if one tries to set the gradient magnitude to 1 in the global Poisson reconstruction.

5.1.2. Quantitative evaluation. As mentioned in [15], there is no reliable criterion to quantitatively evaluate the quality of the result for such reconstruction problems. In our context where only gradient orientations are extracted, it is reasonable to evaluate the reconstruction quality based on the normalized correlation to the input image (which is invariant under affine contrast change). If $u, v : \Omega \rightarrow \mathbb{R}$ are two images, the normalized correlation is

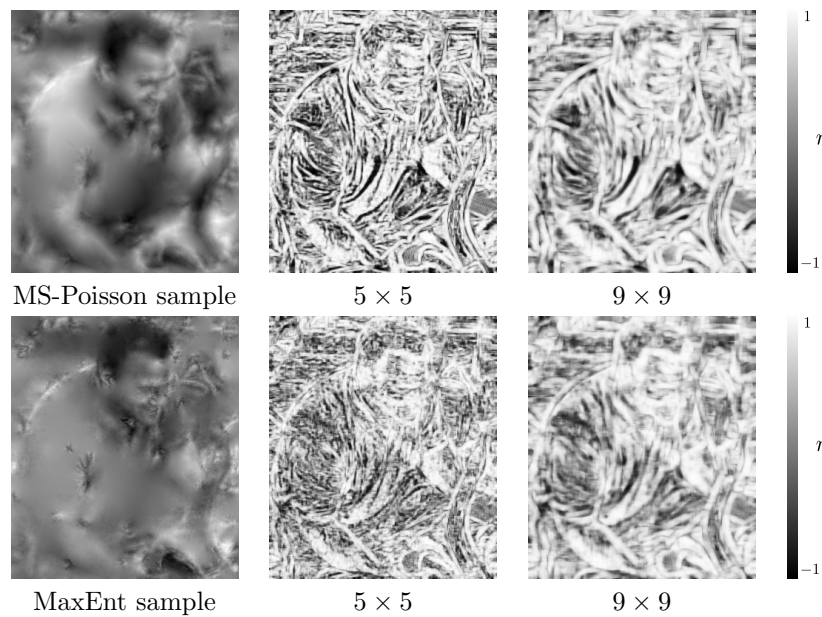


Figure 6. Comparison between MS-Poisson and MaxEnt with local normalized correlations. *In the left column, we display samples of the models MS-Poisson and MaxEnt. In the other columns, we display the local normalized correlation of the sample (first column) with respect to the original image. The local normalized correlations are computed on patch sizes 5×5 and 9×9 , with values in $[-1, 1]$. See the text for comments. (Images are better seen in the electronic version.)*

defined as

$$(53) \quad r(u, v) = \frac{1}{|\Omega|} \sum_{\mathbf{x} \in \Omega} \left(\frac{u(\mathbf{x}) - \bar{u}}{\sigma_u} \right) \left(\frac{v(\mathbf{x}) - \bar{v}}{\sigma_v} \right) \in [-1, 1],$$

where $\bar{u} = \frac{1}{|\Omega|} \sum_{\mathbf{x} \in \Omega} u(\mathbf{x})$ and $\sigma_u^2 = \frac{1}{|\Omega|} \sum_{\mathbf{x} \in \Omega} (u(\mathbf{x}) - \bar{u})^2$. In Figure 4, for each result we have indicated the normalized correlation value r . Surprisingly, the higher correlation values are attained with results linked to the MS-Poisson model (even if it only has access to HOGs computed on a blurred gradient). In addition, the value attained by the samples (or mean) of MS-Poisson is close to the one obtained with the true multiscale HOGs. In contrast, the correlations obtained with the MaxEnt model are lower. This is better explained by the results of Figure 6, in which we display values of local normalized correlations obtained with both models: for each pixel \mathbf{x} we extract patches $p_{\mathbf{x}}(u), p_{\mathbf{x}}(v)$ of compared images u, v and compute the normalized correlation $r(p_{\mathbf{x}}(u), p_{\mathbf{x}}(v))$. On the one hand, the MaxEnt result is everywhere much noisier (because gradient orientations are sampled independently). On the other hand, the regularization involved in MS-Poisson helps to propagate good correlations values in regions located near SIFT subcells. Also, this criterion based on normalized correlation confirms the choice for the regularization parameter $\mu = 50$; see Figure 5.

Another interesting way of performing quantitative evaluation in our context is to compare the HOGs computed in the SIFT subcells to those of the original image. For each subcell, we can compute histograms H_u, H_v of gradient orientations (with 8 bins) for the

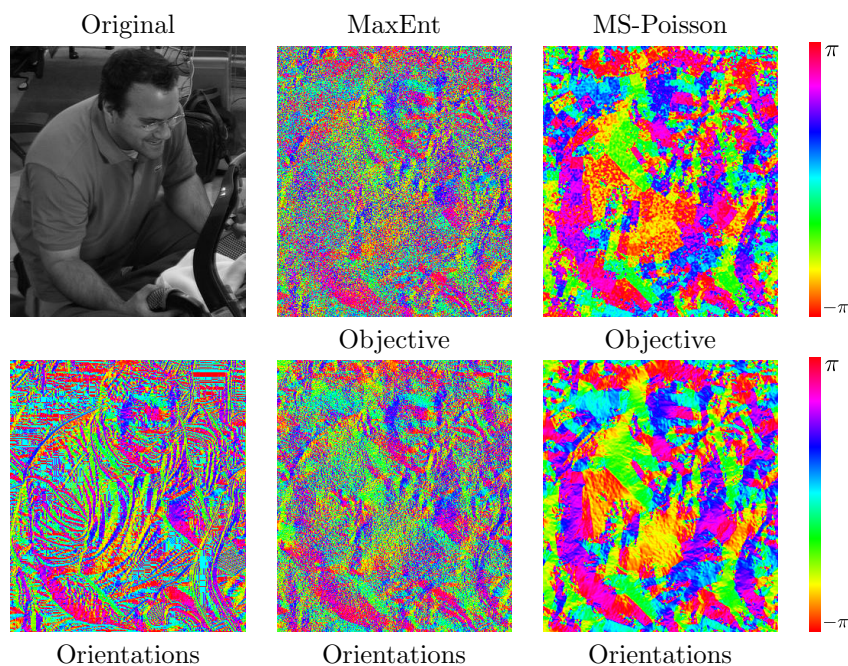


Figure 7. Orientation fields of the original and reconstructions. *In the left column, we display the original image (top) and the corresponding gradient orientations (bottom). In the middle (MaxEnt) and right (MS-Poisson) columns, we display the orientation of the objective vector field (before Poisson reconstruction, top) and the orientation of the resulting gradient field (after Poisson reconstruction, bottom). In the regions that are covered by several SIFT subcells, one can see that the local HOGs are quite well preserved (especially for MaxEnt), even if the orientations are locally shuffled. One can also observe that the Poisson reconstruction step smooths slightly the orientation field. (Images are better seen in the electronic version.)*

images u, v and then compute the total variation distance between these histograms, defined as $\frac{1}{2} \sum_{\ell=1}^8 |H_u(\ell) - H_v(\ell)| \in [0, 1]$. Again, we use gradients at scale 0 when considering the MaxEnt model, and scaled gradients when considering the MS-Poisson model. We can then average the HOG distances obtained for all SIFT subcells, weighted by the number of pixels in each subcell. With this methodology, for the image of Figure 4, we obtain a mean distance around 0.27 for MS-Poisson and 0.16 for MaxEnt. This value is lower for MaxEnt because the model is inherently made to satisfy the HOG constraint. One can better understand these results by examining the orientation fields of both models, as proposed in Figure 7, in particular the effect of the final Poisson reconstruction (keeping in mind that MS-Poisson can also be written with a single objective vector field given by (43)). In this figure, one clearly observes that the objective vector field for MS-Poisson is already very smooth (and certainly too smooth to account for fine local variations in orientation). In contrast the objective vector field for MaxEnt better accounts for the fine variations, but is much noisier, even after the Poisson reconstruction step.

5.1.3. Second order statistics. As we have seen in section 4.3, it is possible to compute the second order statistics of the reconstructed image in each model. In Figure 8 we display the standard deviations of all pixel values in each model. One first remark is that MaxEnt has

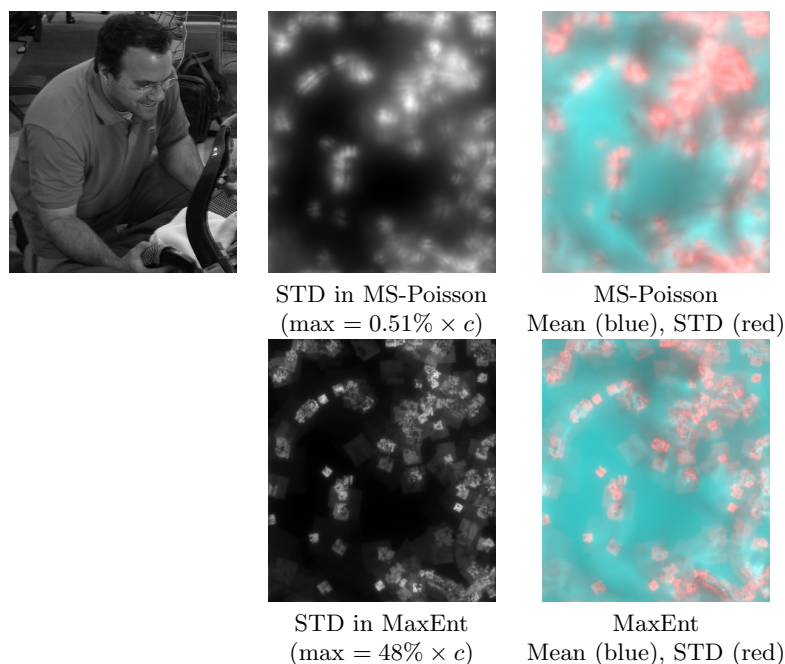


Figure 8. Standard deviations of MS-Poisson and MaxEnt models. *On the top left we display the original image. In the rest of the figure we display the images formed with the standard deviations (STD) of the models MS-Poisson (first row) and MaxEnt (second row). In the second column, we display the raw STD values. In the third column, the red component corresponds to the raw STD values (same as in the second column), and the blue component corresponds to the mean image $m = \mathbb{E}(U)$ of the model (MaxEnt or MS-Poisson). Let us emphasize that for better visualization the images of the second column are renormalized so that the white color corresponds to the indicated maximum value (expressed as a percentage of the empirical standard deviation $c = \sqrt{|\Omega|^{-1} \sum m(\mathbf{x})^2 - (|\Omega|^{-1} \sum m(\mathbf{x}))^2}$ of the mean image m). These results clearly indicate that the MS-Poisson model is much more concentrated around its expectation than MaxEnt. (Images are better seen in the electronic version.)*

in general a much larger variance than MS-Poisson, which can be explained by the fact that the output of MS-Poisson is in some sense a weighted average of many local reconstructions. Also it is interesting to see that the image regions with larger variance are located in the SIFT subcells which contain sharp geometric details. That being said, the variance of both of these models is relatively small compared to the global range of the mean image, which indicates that both models have quite small variations around the mean.

5.1.4. Discard boundary keypoints. Let us emphasize that in our experiments, we used all the keypoints computed by the SIFT methods and did not discard keypoints located near the image boundaries. The positions of the corresponding local extrema in the normalized scale-space are indeed highly dependent on the boundary conditions used to compute the scale-space. This explains why SIFT keypoints near the image boundaries are often discarded for particular applications, e.g., image matching. In our reconstruction problem, there is no reason to discard such keypoints, and we use the information available in SIFT subcells as soon as they intersect the image domain (if the SIFT subcell is not entirely contained in the domain, we consider only the pixels in the intersection of the subcell and the domain). But

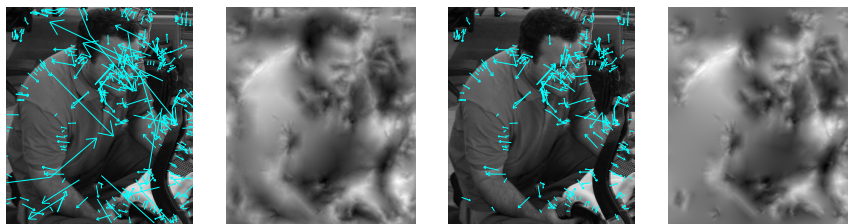


Figure 9. Discard keypoints near image boundary. *In this figure, we examine the effect of discarding keypoints whose associated SIFT cell is not entirely contained in the image domain. The displayed reconstructions are samples of the MS-Poisson model.*

still, it is clear that for some images the reconstruction will be quite different when discarding those keypoints. For example, in the case of Figure 9, if boundary keypoints are discarded, then several parts of the man’s body are not as properly retrieved in the reconstruction, thus affecting the semantic understanding of the image.

5.1.5. Matching keypoints between the original and reconstructed images. Finally, it is interesting to compare the keypoints computed on the original image and the ones computed on several samples of the models. As one can see in Figure 10, we get back similar keypoints in many regions of the image, but still with some variations in position, scale, and orientation. In particular, we observe variations when taking different samples of the model (sometimes some keypoints associated with low contrast regions may even disappear). Notice also that we get back fewer keypoints in the MS-Poisson model; indeed, since it is more regular we lose some extrema in the scale-space. In addition, the regularization tends to change the scale of the structures, and thus the scales of the keypoints is often larger than in the original image.

In order to give a more quantitative evaluation of the variations of the keypoints over different samples of the model, it is possible to use the matching algorithm available with the online implementation [53] (we used the proposed default parameters). This algorithm follows the matching method proposed in [32], which essentially pairs SIFT keypoints by thresholding the ratio between the distances to the first and second nearest neighbors (computed with the ℓ^2 -distance between SIFT descriptors). First we can comment on what happens when matching two different samples of the same model. For the MS-Poisson model, when matching the two samples shown in Figure 10, among the 206 keypoints found on the first image (resp., 211 on the second image), 150 keypoints are matched. The mean spatial distance (resp., mean scale variation, mean angle variation) between matched keypoints is about 0.54 (resp., 0.15, 0.050). Similar numbers can be given for the MaxEnt model, but in this case much fewer keypoints are correctly matched: of the 452 keypoints found on the first image (resp., 458 on the second image), only 184 are matched. This reflects again the larger variance of the MaxEnt model.

More interestingly, we can try to match the SIFT keypoints between the original image and the reconstructions. Unfortunately, only a few SIFT points are properly matched this way: among the 477 keypoints found in the original image, around 10 keypoints are properly matched in samples of the MS-Poisson model, and no keypoints are matched when comparing to a sample of MaxEnt. This shows that even if these models are able to recover gradient

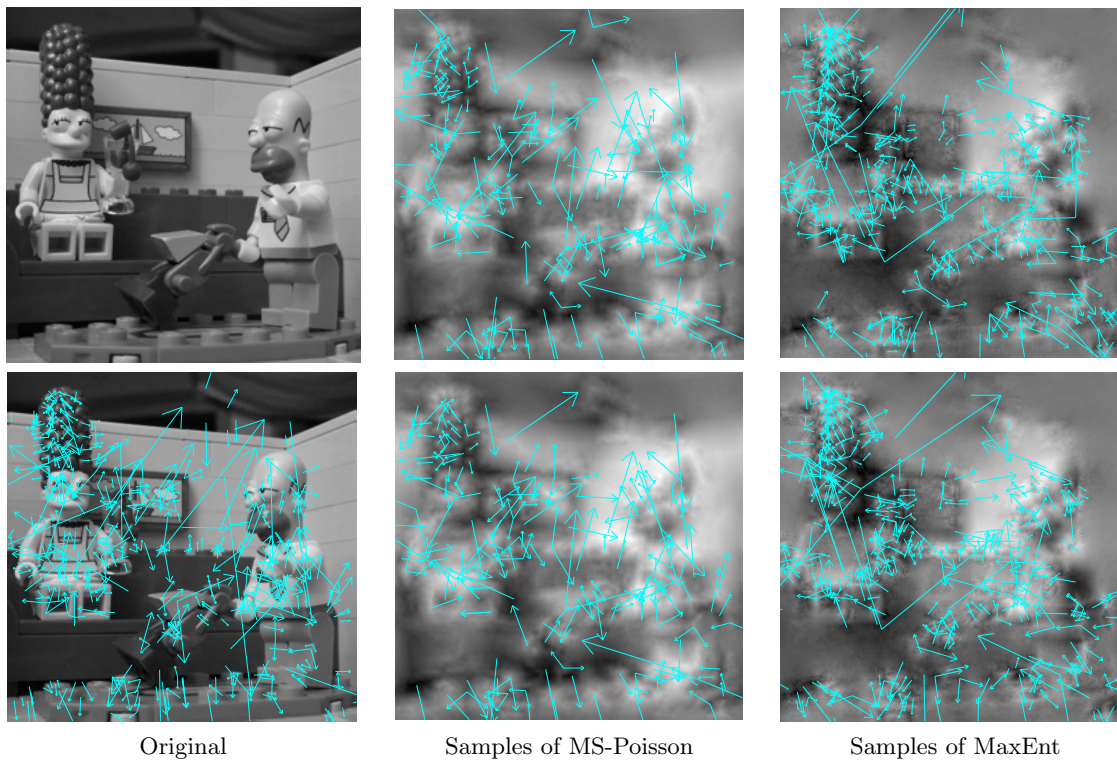


Figure 10. Keypoints after reconstruction. *In the first column, we display an original image (courtesy of J. Delon) and the same image with its SIFT keypoints. In the second column, we display two samples of the MS-Poisson model. In the third column, we display two samples of the MaxEnt model. We display the keypoints associated to these images as superimposed blue arrows. Notice that several keypoints are retrieved after reconstruction, with still some variations in position and orientation. Notice also that we observe some variations in the keypoints associated to different samples of these models. See the text for additional comments. (Images are better seen in the electronic version.)*

orientations in a somehow blurry manner, this is not sufficient to precisely get back the content of SIFT descriptors. By the way, the fact that only 75% (resp., 50%) of the keypoints are matched between two samples of MS-Poisson (resp., MaxEnt) illustrates the sensitivity of the SIFT descriptors to small random perturbations.

5.2. Reconstruction from true SIFT descriptors. The two models MS-Poisson and MaxEnt are designed to propose stochastic reconstructions of an image based on simplified SIFT descriptors, that is, multiscale HOGs extracted around the SIFT keypoints. But it is also possible to test these reconstruction models with the true SIFT descriptors. For that, for each keypoint, we still consider the location, scale, and principal orientation, but, following the discussion of section 2.2, starting from the normalized feature vector $(f_k) \in \mathbb{R}^{128}$, we improperly build target histograms for the 16 corresponding SIFT subcells: for each $p \in \{1, \dots, 16\}$, to the corresponding p th subcell s_j , we associate the discrete histogram

$$(54) \quad \tilde{H}_{j,\ell} = \frac{f_{16(p-1)+\ell}}{\sum_{\ell'=1}^8 f_{16(p-1)+\ell'}} \quad (1 \leq \ell \leq 8).$$

We can thus sample the MS-Poisson model using the $(\tilde{H}_{j,\ell})$ values as a substitute for the extracted multiscale HOG $(H_{j,\ell})$.

In Figure 11, we display several reconstruction results obtained with the MS-Poisson model based on the multiscale HOGs or the true SIFT descriptors. As could be expected, the reconstruction results obtained with the true SIFT descriptors are not as good as those obtained from multiscale HOGs; in particular many fine-scale structures are lost, and the shape of small objects is not recovered in a coherent way (see, for example, the wings in the butterfly image). However, large-scale structures of the image are still retrieved quite properly, which often suffices to understand the semantic content of the image.

In order to get sharper results, we should adapt the reconstruction models to account for the normalizations applied in the original SIFT method. It appears quite straightforward to adapt the models to histograms computed with linear votes (instead of binary votes). However, it seems much more difficult to cope with the final normalization and thresholding (see (2)), which dramatically reduce the quantity of information. Also, in the true SIFT descriptors, the pixels vote for orientation values with a weight that is proportional to the gradient magnitude. This explains why it is difficult to retrieve the local HOG from the SIFT descriptors in the absence of any information about the local gradient magnitude.

5.3. Comparison with previous works. In this subsection, we propose to compare our reconstruction models with those obtained by the methods of Weinzaepfel, Jégou, and Pérez [64] and Dosovitskiy and Brox [19]; see Figure 12. One important difference between these two other approaches and ours is that our method relies only on the content provided in the SIFT subcells, while these methods exploit an external database either to copy local information from patches with similar SIFT descriptors (as in [64]) or to build an up-convolutional neural network for reconstruction (as in [19]). Thus our work has no intention to outperform these methods in terms of visual quality of reconstruction (in particular, our method has absolutely no possibility of recovering the color information). Notice that we cannot compare to the method of [35], which is adapted to “dense SIFT” (i.e., SIFT descriptors computed on a dense set of patches) and not “sparse SIFT” (i.e., SIFT descriptors computed around the keypoints).

There are also minor differences in the extracted information because both of these works do not rely on the original implementation of the SIFT method. The method of [64] actually uses “elliptic” interest regions (extracted using the Hessian-affine method by [41]) in which normalized multiscale HOG are computed (in the same way as in the original SIFT method). In contrast, Dosovitskiy and Brox use circular keypoints and descriptors that are computed with the VLFeat library [61]. But in order to apply an up-convolutional neural network to these features, they need to derive a grid-based representation of these features: the image is divided into 4×4 cells, and each cell containing a keypoint is associated with the corresponding oriented keypoint and feature vector. If there is no keypoint, then they associate the zero vector, and if there are several keypoints, they randomly choose one of them (see the details in [19, section III]).

One advantage of the MS-Poisson model, compared to the result of [64], is that it is defined through the minimization of the global MS-Poisson energy (37). Therefore, it produces images that are globally coherent while respecting as much as possible the local constraints given by

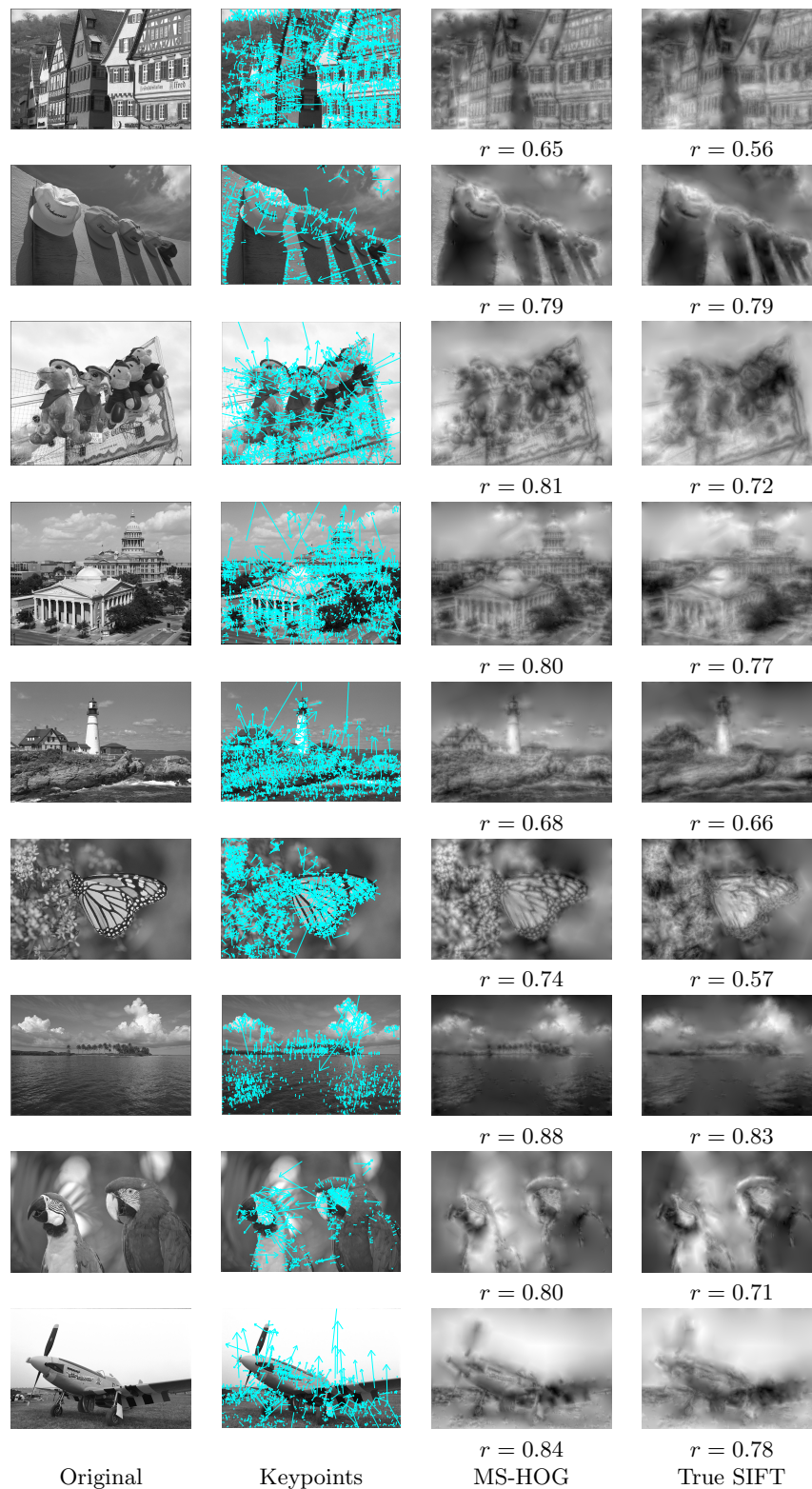


Figure 11. Reconstruction results from multiscale HOG or SIFT descriptors with images from the Live database [56]. For each row, from left to right, we display an original image, the same image with superimposed SIFT keypoints, a sample of the MS-Poisson model obtained from multiscale HOG, and a sample of the MS-Poisson model obtained from the true SIFT descriptors. Notice that the reconstruction from true SIFT descriptors is less sharp but still recovers many geometric structures of the initial image.

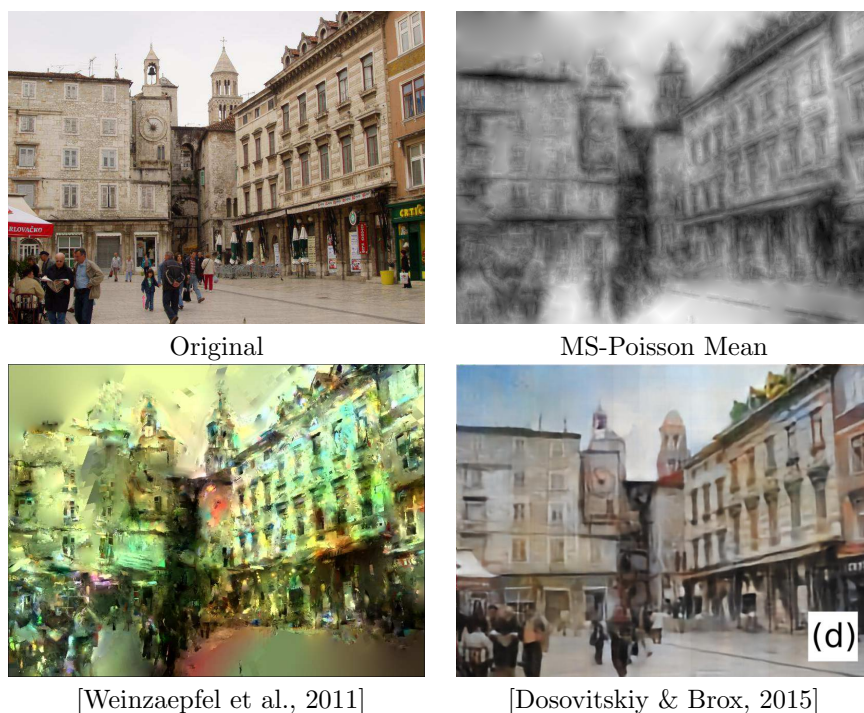


Figure 12. Comparison for SIFT reconstruction. *In the first row we display the original image and the reconstruction results obtained as the expectation of the MS-Poisson model computed on the true SIFT descriptors (see section 5.2). In the second row we display the results obtained with the methods of [64] and [19]. Notice that the MS-Poisson model provides images that are blurrier but also more globally coherent than the ones obtained by the method of [64]. However, this model does not compete with [19] in terms of restitution and visual quality since it does not rely on any external information.*

the multiscale HOGs. In contrast, the result of [64] is clearly affected by stitching artifacts which are inherent to their reconstruction method. On the other hand, their method is able to copy pieces of clean patches so that their reconstruction looks locally sharper (but also noisier).

However, the reconstructed images obtained in [19] are both globally coherent and quite sharp. Indeed, our method does not rely on an external database, so it cannot compete with that of [19], and in particular it cannot get back information which is completely lost in the SIFT descriptors (global contrast, or also color information).

5.4. Reconstruction with other keypoints. In this subsection, we question the very definition of the SIFT keypoints in terms of synthesis, similarly to what was done in [49]. Indeed, one can wonder if selecting the local extrema of $(\mathbf{x}, \sigma) \mapsto \sigma^2 \Delta g_\sigma * u(\mathbf{x})$ is the best possible choice for points of interest in order to extract relevant information for synthesis.

For that, we propose a comparison with two other sets of keypoints extracted in a very different way. The first choice (“Min-Rec-Error”) is driven by the following intuition: using the Taylor formula around a point \mathbf{x} , one can write when $\sigma \rightarrow 0$ that

$$(55) \quad \int u(\mathbf{x} + \mathbf{z}) g_\sigma(\mathbf{z}) d\mathbf{z} - u(\mathbf{x}) = \sigma^2 \Delta u(\mathbf{x}) + o(\sigma^2).$$

Therefore, near the positions \mathbf{x} where $\Delta u(\mathbf{x})$ is close to zero, one can approximately recover $u(\mathbf{x})$ by averaging neighboring values. In this sense, it seems relevant to extract more information at the points where the average reconstruction fails, in particular at the maxima of $|\Delta u|$.

But one could also directly work with the reconstruction error; we thus propose to extract local maxima of the function

$$(56) \quad (\mathbf{x}, \sigma) \mapsto |g_\sigma * u(\mathbf{x}) - u(\mathbf{x})|.$$

In our implementation, we detect these maxima on a discretized scale-space with 30 scales $s = 2^{r/6}$, $0 \leq r < 30$. Also, in order to draw a comparison with a fixed number of keypoints, we only keep points having an “edgeness” value below a threshold. As in the original SIFT method, the edgeness measure is obtained as the ratio $\frac{\text{Tr}(H)^2}{\det H}$ of the principal curvatures, where H is the Hessian of the smoothed image $g_2 * u$. The threshold is adapted in order to get the same number n_{kp} of keypoints as that provided by the SIFT method.

The second and third choices (“Random-unif” and “Random-grad”) consist in selecting keypoints in a random manner. More precisely, for the choice Random-unif, we independently sample n_{kp} keypoints by choosing uniformly a position \mathbf{x} in the image domain, a uniform orientation $\alpha \in \mathbb{T}$, and a scale by sampling an exponential distribution whose parameter is adjusted so that the expectation is the same as the mean scale of the usual SIFT keypoints. Modeling by the exponential distribution is empirically justified by the fact that the distribution of scales of SIFT keypoints is concentrated in the fine scales. For the choice Random-grad, we do the same except that the positions are randomly drawn using a probability distribution which is proportional to the gradient magnitude of the smoothed image $g_2 * u$.

For these new sets of keypoints, we computed the average image of the MS-Poisson model. The results are displayed in Figure 13. They clearly indicate that the usual SIFT keypoints lead to a reconstruction that is visually better than the others. The main problem of the Min-Rec-Error keypoints is that they do not extract enough small-scale information: for the examples shown in Figure 13 the average scale of these keypoints is approximately twice as large as that of the SIFT keypoints. In addition, for both Min-Rec-Error and random keypoints, the spatial locations are not concentrated around geometric details, as can be the case with the SIFT keypoints. The comparison with Random-grad is particularly interesting; indeed the reconstruction with Random-grad keypoints is slightly better than that with Random-unif keypoints, but still it fails to recover fine details. The main problem of the Random-grad approach is that it is not contrast invariant, and thus it favors points with strong gradients in uniform regions over points in salient regions with low contrast. Thus, the usual definition of SIFT keypoints (and in particular the thresholding steps) is confirmed to be a relevant choice for extracting visual information near salient structures, from both analysis- and synthesis-based perspectives.

6. Conclusion. In this paper we proposed two stochastic models (MaxEnt, resp., MS-Poisson) for reconstructing an image based only on the information contained in the (monoscale, resp., multiscale) local HOGs computed in the SIFT subcells. With both models we get back images which are close to the original in terms of semantic content. This is still true if we compute the reconstructions based on the true SIFT descriptors. One benefit of these mod-

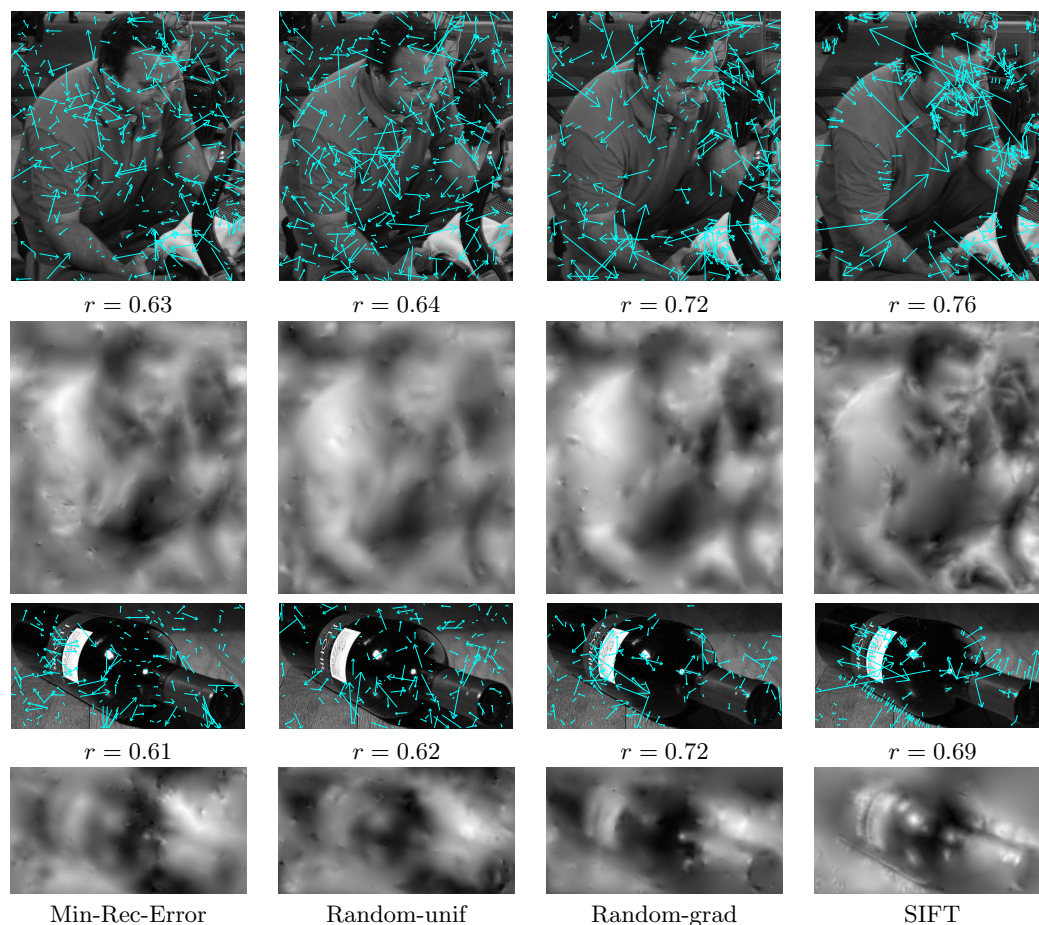


Figure 13. Reconstruction with other keypoints. The first column (“Minimum reconstruction error”) corresponds to the keypoints obtained as local minima of (56). The second (“Random-unif”) and third (“Random-grad”) columns correspond to the randomly selected keypoints. The last column corresponds to the standard SIFT keypoints. The original images are displayed in Figures 3 and 4. Above each reconstruction we indicate the value of the normalized correlation to the original image. See the text in section 5.4 for the precise definition of these sets of keypoints, and additional comments.

els over competing approaches is that they do not rely on any external image database, and besides the convolutive expressions found in this paper allow one to compute the statistics of the corresponding output random fields (e.g., local variance).

However, several questions raised by this work remain open. First it would be interesting to consider generalizations of the MS-Poisson model with different image priors, i.e., adopt other regularization terms in the functional. It is likely that solving the corresponding optimization problem may require an iterative procedure, but on the other hand the solutions may exhibit cleaner geometric structures which are better extrapolated outside the SIFT subcells. Also, there is more to discuss about the optimality of keypoints with respect to the quality of reconstructed images. In particular, here we adopted one unique reconstruction strategy in order to compare different sets of keypoints. But it seems possible to optimize both the sets

of keypoints and the reconstruction strategy in order to maximize a criterion linked to the proximity of the reconstruction to the input original image. This could be thought of as a kind of auto-encoding procedure in which the encoder is constrained to have a very particular form (that is, keypoint extractor).

REFERENCES

- [1] T. AHONEN, A. HADID, AND M. PIETIKAINEN, *Face description with local binary patterns: Application to face recognition*, IEEE Trans. Pattern Anal. Mach. Intell., 28 (2006), pp. 2037–2041.
- [2] A. ALAHI, R. ORTIZ, AND P. VANDERGHEYNST, *FREAK: Fast retina keypoint*, in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, 2012, pp. 510–517.
- [3] B. ALLEN AND M. KON, *The Marr Conjecture and Uniqueness of Wavelet Transforms*, preprint, <https://arxiv.org/abs/1401.0542>, 2015.
- [4] B. ALLEN AND M. KON, *Unique recovery from edge information*, in 2015 International Conference on Sampling Theory and Applications (SampTA), IEEE, 2015, pp. 312–316.
- [5] F. ATTNEAVE, *Some informational aspects of visual perception*, Psychological Rev., 61 (1954), pp. 183–193.
- [6] S. BATTIATO, G. GALLO, G. PUGLISI, AND S. SCCELLATO, *SIFT features tracking for video stabilization*, in 14th International Conference on Image Analysis and Processing (ICIAP 2007), IEEE, 2007, pp. 825–830.
- [7] H. BAY, A. ESS, T. TUYTELAARS, AND L. VAN GOOL, *Speeded-up robust features (SURF)*, Comput. Vis. Image Understanding, 110 (2008), pp. 346–359.
- [8] M. BLACK AND A. JEPSON, *Eigentracking: Robust matching and tracking of articulated objects using a view-based representation*, Int. J. Comput. Vis., 26 (1998), pp. 63–84.
- [9] Y.-L. BOUREAU, F. BACH, Y. LECUN, AND J. PONCE, *Learning mid-level features for recognition*, in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, 2010, pp. 2559–2566.
- [10] Y.-L. BOUREAU, N. LE ROUX, F. BACH, J. PONCE, AND Y. LECUN, *Ask the locals: Multi-way local pooling for image recognition*, in 2011 IEEE International Conference on Computer Vision (ICCV 2011), IEEE, 2011, pp. 2651–2658.
- [11] G. CSURKA, C. DANCE, L. FAN, J. WILLAMOWSKI, AND C. BRAY, *Visual categorization with bags of keypoints*, in Workshop on Statistical Learning in Computer Vision, ECCV, 2004.
- [12] S. CURTIS AND A. OPPENHEIM, *Reconstruction of multidimensional signals from zero crossings*, J. Opt. Soc. Amer. A, 4 (1987), pp. 221–231, <https://doi.org/10.1364/JOSAA.4.000221>.
- [13] S. CURTIS, S. SHITZ, AND A. OPPENHEIM, *Reconstruction of nonperiodic two-dimensional signals from zero crossings*, IEEE Trans. Acoust. Speech Signal Process., 35 (1987), pp. 890–893.
- [14] N. DALAL AND B. TRIGGS, *Histograms of oriented gradients for human detection*, in 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05), Vol. 1, IEEE, 2005, pp. 886–893.
- [15] E. D'ANGELO, L. JACQUES, A. ALAHI, AND P. VANDERGHEYNST, *From bits to images: Inversion of local binary descriptors*, IEEE Trans. Pattern Anal. Mach. Intell., 36 (2014), pp. 874–887.
- [16] A. DESOLNEUX, *When the a contrario approach becomes generative*, Int. J. Comput. Vis., 116 (2016), pp. 46–65.
- [17] A. DESOLNEUX AND A. LECLAIRE, *Stochastic image reconstruction from local histograms of gradient orientation*, in Proceedings of the Sixth International Conference on Scale Space and Variational Methods in Computer Vision (SSVM), Lecture Notes in Comput. Sci. 10302, Springer, 2017, pp. 133–145.
- [18] A. DESOLNEUX, L. MOISAN, AND J. MOREL, *From Gestalt Theory to Image Analysis: A Probabilistic Approach*, Interdiscip. Appl. Math. 34, Springer Science & Business Media, 2007.
- [19] A. DOSOVITSKIY AND T. BROX, *Inverting Visual Representations with Convolutional Networks*, preprint, <https://arxiv.org/abs/1506.02753>, 2015.
- [20] J. H. ELDER AND S. W. ZUCKER, *Scale space localization, blur, and contour-based image coding*, in Proceedings of the 1996 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'96), IEEE, 1996, pp. 27–34.

- [21] O. FAUGERAS, *Three-Dimensional Computer Vision: A Geometric Viewpoint*, MIT Press, 1993.
- [22] P. FELZENSZWALB, R. GIRSHICK, D. MCALLESTER, AND D. RAMANAN, *Object detection with discriminatively trained part-based models*, IEEE Trans. Pattern Anal. Mach. Intell., 32 (2010), pp. 1627–1645.
- [23] C. HARRIS AND M. STEPHENS, *A combined corner and edge detector*, in Proc. 4th Alvey Vision Conference, Citeseer, 1988, pp. 147–151.
- [24] R. HUMMEL AND R. MONIOT, *Reconstructions from zero crossings in scale space*, IEEE Trans. Acoust. Speech Signal Process., 37 (1989), pp. 2111–2130.
- [25] F. JUEFEI-XU AND M. SAVVIDES, *Learning to invert local binary patterns*, in the 27th British Machine Vision Conference (BMVC), 2016.
- [26] H. KATO AND T. HARADA, *Image reconstruction from bag-of-visual-words*, in Proceedings of the 2014 IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2014, pp. 955–962.
- [27] A. KRIZHEVSKY, I. SUTSKEVER, AND G. HINTON, *ImageNet classification with deep convolutional neural networks*, in Advances in Neural Information Processing Systems, 2012, pp. 1097–1105.
- [28] S. LAZEBNIK, C. SCHMID, AND J. PONCE, *Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories*, in IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Vol. 2, IEEE, 2006, pp. 2169–2178.
- [29] S. LEUTENEGGER, M. CHLI, AND R. Y. SIEGWART, *BRISK: Binary robust invariant scalable keypoints*, in IEEE International Conference on Computer Vision (ICCV), IEEE, 2011, pp. 2548–2555.
- [30] T. LINDBERG, *Feature detection with automatic scale selection*, Int. J. Comput. Vis., 30 (1998), pp. 79–116.
- [31] T. LINDBERG, *Image matching using generalized scale-space interest points*, J. Math. Imaging Vision, 52 (2015), pp. 3–36.
- [32] D. LOWE, *Distinctive image features from scale-invariant keypoints*, Int. J. Comput. Vis., 60 (2004), pp. 91–110.
- [33] Y. LU, S. ZHU, AND Y. N. WU, *Learning FRAME Models Using CNN Filters for Knowledge Visualization*, preprint, <https://arxiv.org/abs/1509.08379>, 2015.
- [34] A. MAHENDRAN AND A. VEDALDI, *Understanding deep image representations by inverting them*, in Proceedings of the 2015 IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2015, pp. 5188–5196.
- [35] A. MAHENDRAN AND A. VEDALDI, *Visualizing deep convolutional neural networks using natural pre-images*, Int. J. Comput. Vis., 120 (2016), pp. 233–255.
- [36] E. MAIR, G. D. HAGER, D. BURSCHKA, M. SUPPA, AND G. HIRZINGER, *Adaptive and generic corner detection based on the accelerated segment test*, in European Conference on Computer Vision, Springer, 2010, pp. 183–196.
- [37] S. MALLAT AND S. ZHONG, *Characterization of signals from multiscale edges*, IEEE Trans. Pattern Anal. Mach. Intell., 14 (1992), pp. 710–732.
- [38] D. MARR, *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information*, W.H. Freeman and Company, 1982.
- [39] D. MARR AND E. HILDRETH, *Theory of edge detection*, Proc. R. Soc. Lond. Ser. B Biol. Sci., 207 (1980), pp. 187–217.
- [40] Y. MEYER, *Wavelets: Algorithms and Applications*, SIAM, 1993.
- [41] K. MIKOLAJCZYK AND C. SCHMID, *Scale & affine invariant interest point detectors*, Int. J. Comput. Vis., 60 (2004), pp. 63–86.
- [42] K. MIKOLAJCZYK AND C. SCHMID, *A performance evaluation of local descriptors*, IEEE Trans. Pattern Anal. Mach. Intell., 27 (2005), pp. 1615–1630.
- [43] J.-M. MOREL, A. PETRO, AND C. SBERT, *Fourier implementation of Poisson image editing*, Pattern Recognition Lett., 33 (2012), pp. 342–348.
- [44] J.-M. MOREL AND G. YU, *ASIFT: A new framework for fully affine invariant image comparison*, SIAM J. Imaging Sci., 2 (2009), pp. 438–469, <https://doi.org/10.1137/080732730>.
- [45] J.-M. MOREL AND G. YU, *Is SIFT scale invariant?*, Inverse Probl. Imaging, 5 (2011), pp. 115–136.
- [46] D. MUMFORD AND A. DESOLNEUX, *Pattern Theory: The Stochastic Analysis of Real-World Signals*, A K Peters/CRC Press, 2010.
- [47] P. MUSÉ, F. SUR, F. CAO, Y. GOUSSEAU, AND J.-M. MOREL, *An a contrario decision method for shape element recognition*, Int. J. Comput. Vis., 69 (2006), pp. 295–315.

- [48] Y. NESTEROV, *Introductory Lectures on Convex Optimization: A Basic Course*, Appl. Optim. 87, Springer, 2004.
- [49] M. NIELSEN AND M. LILLHOLM, *What do features tell about images?*, in *Scale-Space*, Vol. 1, Springer, 2001, pp. 39–50.
- [50] T. OJALA, M. PIETIKÄINEN, AND T. MÄENPÄÄ, *Multiresolution gray-scale and rotation invariant texture classification with local binary patterns*, IEEE Trans. Pattern Anal. Mach. Intell., 24 (2002), pp. 971–987.
- [51] P. PÉREZ, M. GANGNET, AND A. BLAKE, *Poisson image editing*, in *ACM SIGGRAPH 2003 Papers*, SIGGRAPH '03, 2003, pp. 313–318, <https://doi.org/10.1145/1201775.882269>.
- [52] J. PHILBIN, O. CHUM, M. ISARD, J. SIVIC, AND A. ZISSERMAN, *Object retrieval with large vocabularies and fast spatial matching*, in *Proceedings of the 2007 IEEE Conference on Computer Vision and Pattern Recognition*, IEEE, 2007, pp. 1–8.
- [53] I. REY OTERO AND M. DELBRACIO, *Anatomy of the SIFT method*, Image Process. On Line, 4 (2014), pp. 370–396, <https://doi.org/10.5201/ipol.2014.82>.
- [54] E. ROSTEN AND T. DRUMMOND, *Machine learning for high-speed corner detection*, in *Computer Vision—ECCV 2006*, Springer, 2006, pp. 430–443.
- [55] J. SANZ AND T. HUANG, *Theorems and Experiments on Image Reconstruction from Zero Crossings*, IBM Almaden Research Center, 1987.
- [56] H. SHEIKH, Z. WANG, L. CORMACK, AND A. BOVIK, *Live Image Quality Assessment Database*, release 2 (2005), <http://live.ece.utexas.edu/research/quality/subjective.htm>.
- [57] K. SIMONYAN AND A. ZISSERMAN, *Very deep convolutional networks for large-scale image recognition*, in *Proceedings of the International Conference on Learning Representations*, 2014.
- [58] J. SIVIC AND A. ZISSERMAN, *Video Google: A text retrieval approach to object matching in videos*, in *Proceedings of the 9th IEEE International Conference on Computer Vision*, IEEE, 2003, pp. 1470–1477.
- [59] T. TUYTELAARS AND K. MIKOLAJCZYK, *Local invariant feature detectors: A survey*, Found. Trends Comput. Graphics Vis., 3 (2008), pp. 177–280.
- [60] T. TUYTELAARS AND L. VAN GOOL, *Matching widely separated views based on affine invariant regions*, Int. J. Comput. Vis., 59 (2004), pp. 61–85.
- [61] A. VEDALDI AND B. FULKERSON, *VLFeat: An open and portable library of computer vision algorithms*, in *Proceedings of the 18th ACM International Conference on Multimedia*, ACM, 2010, pp. 1469–1472.
- [62] C. VONDRICK, A. KHOSLA, T. MALISIEWICZ, AND A. TORRALBA, *HOGgles: Visualizing object detection features*, in *Proceedings of the 2013 IEEE International Conference on Computer Vision*, 2013, pp. 1–8.
- [63] C. WALLRAVEN, B. CAPUTO, AND A. GRAF, *Recognition with local features: The kernel recipe*, in *Proceedings of the 9th IEEE International Conference on Computer Vision*, IEEE, 2003, pp. 257–264.
- [64] P. WEINZAEPFEL, H. JÉGOU, AND P. PÉREZ, *Reconstructing an image from its local descriptors*, in *Proceedings of the IEEE CVPR*, 2011, pp. 337–344.
- [65] J. YANG, D. SCHONFELD, AND M. MOHAMED, *Robust video stabilization based on particle filter tracking of projected camera motion*, IEEE Trans. Circuits Systems Video Technol., 19 (2009), pp. 945–954.
- [66] A. YILMAZ, O. JAVED, AND M. SHAH, *Object tracking: A survey*, ACM Comput. Surveys, 38 (2006), 13.
- [67] M. ZEILER AND R. FERGUS, *Visualizing and understanding convolutional networks*, in *Computer Vision—ECCV 2014*, Springer, 2014, pp. 818–833.
- [68] J. ZHANG, M. MARSZALEK, S. LAZEBNIK, AND C. SCHMID, *Local features and kernels for classification of texture and object categories: A comprehensive study*, Int. J. Comput. Vis., 73 (2007), pp. 213–238.
- [69] S. ZHU, Y. WU, AND D. MUMFORD, *Filters, random fields and maximum entropy (FRAME): Towards a unified theory for texture modeling*, Int. J. Comput. Vis., 27 (1998), pp. 107–126.