



HAL
open science

Stochastic Image Models from SIFT-like descriptors

Agnès Desolneux, Arthur Leclaire

► **To cite this version:**

Agnès Desolneux, Arthur Leclaire. Stochastic Image Models from SIFT-like descriptors. SIAM Journal on Imaging Sciences, 2018, 10.1137/18M116592X . hal-01692139

HAL Id: hal-01692139

<https://hal.science/hal-01692139v1>

Submitted on 24 Jan 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Stochastic Image Models from SIFT-like descriptors

A. Desolneux, A. Leclaire

CMLA, ENS Cachan, CNRS, Université Paris-Saclay, 94235 Cachan, France.

January 24, 2018

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 these models admit convolutive expressions which allow 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.

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 [40, 39] 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 to extract 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, 33, 9, 10, 26], image matching and rectification [33, 60, 32], object detection and tracking [8, 58, 66, 53], video stabilization [6, 65], image classification [29, 68, 28], etc. 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 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).

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

Here we will only mention a few famous local descriptors, and we refer to [43, 59, 45, 32] 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 [42]. The same authors also proposed in [42] 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 [60] 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 [55] the “features from accelerated segment test” (FAST) which is a corner detector accelerated by a machine learning technique. This approach has been further fastened by Mair et al. [37] using optimal decision trees, thus obtaining an “adaptive and generic accelerated segment test” (AGAST). Musé et al. proposed in [48] to extract shapes from the image level lines, and to process them in order to get an affine invariant representation.

In parallel of 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 (LBP) defined by Ojala et al. [51] which extracts signs of differences of image values on pixels located on a circular neighborhood of a keypoint. The LBP were originally designed for texture description but can also be used for face detection [1]. In [33], 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 [45] 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 be also represented as visual words [58, 11] which are obtained as cluster points in a feature space. Later, some authors proposed to describe a patch using local binary descriptors (LBD), which extracts the signs of differences between Gaussian measurements taken at different locations. Using different ways of selecting these locations leads to the methods BRISK [30] (binary robust invariant scalable keypoints) or FREAK [2] (fast retina keypoint). All of these descriptors have quite different invariance properties (evaluated either in a theoretical or experimental framework).

Long before the design of these image descriptors, the question of a minimal representation of an image was thoroughly studied, mainly for compression purpose. Through the concept of *raw primal sketch*, Marr [39] suggested that the human visual system processes images by retaining essentially the lines of zero-crossing of the Laplacian at several scales. This leads to the conjecture that an image is uniquely defined by these zero-crossing lines, a conjecture that was later precised by Mallat [38] using wavelet modulus maxima. Both these conjectures were proved wrong by Meyer [41] but still, algorithms for approximate reconstruction were proposed by Hummel and Moniot [24] for zero-crossings and by Mallat and Zhong [38] for the case of wavelet modulus maxima. Besides, unique characterization can be shown to be true under some additional hypotheses [12, 13, 56, 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 [50] consider the problem of variational reconstruction from linear measurements; in addition to the minimum variance reconstruction (given by the pseudo-inverse of the measurements matrix),

they propose 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 discuss the problem of extracting a subset of linear measurements which leads to the best reconstruction and empirically compare three different strategies for that purpose.

More recently, motivated by privacy issues (since the descriptors may be transmitted on an unsecured network), Weinzaepfel et al. [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 similar description is looked for in the database, and all the patches are stitched together with Poisson image editing [52]. Vondrick et al. [62] address reconstruction from dense HOGs by relying on a paired dictionary representation of HOGs and patches. Also, d’Angelo et al. [15] address 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 [27] formulate reconstruction from bag of visual words as a problem of quadratic assignment. Finally, Juefei-Xu and Savvides [25] propose to invert 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 [28, 67] has urged the need of 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 to build a deconvolution network that allows 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 [35, 36] proposed to search for a pre-image 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 non-linear so that the resulting optimization problem is not convex; so 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] suggest to learn 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 both HOG, SIFT and AlexNet [28] representations. Notice that the inversion/visualization techniques of [67, 19] exploit an external database while the one of [35, 36] 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 pre-images 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. One way to address this problem is to consider the information contained in the descriptors as a statistical measurement, and to sample from the maximum entropy model that complies with these statistical constraints. Such maximal entropy models were considered by Zhu, Wu and Mumford in [69, 47] for texture modelling based on responses to an automatically selected subset of filters chosen in a filter bank. This approach has been recently extended by Lu, Zhu and Wu to responses to a pre-trained neural network [34]. 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 proposes explicit computations of maximal entropy image models that lead to the same detections (in average).

In the present paper, we propose two stochastic models that complies 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 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 to coherently merge information given at several scales. Several experiments show that both these 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 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 to comply with multiscale HOGs taken in the SIFT subcells; the corresponding H^1 -regularized multiscale Poisson problem is explicitly solved. Finally, in Section 5 we display several reconstruction results obtained with both models (applied with simplified SIFT, or also the true SIFT), 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 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 open some perspectives for future research. A preliminary version of this work was published as a conference paper in [17].

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 later used for the reconstruction in the next sections.

2.1 Gaussian Scale-Space and Keypoints

Following [31], 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 \mathcal{C}^∞ function on $\mathbb{R}^2 \times (0, \infty)$).

Then we consider the scale-normalized Laplacian of Gaussian $\sigma^2\Delta g_\sigma$. The PDE satisfied by k_t gives after 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 2.1 ([31]). *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 [46] 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 [33], the scale-normalized LoG is approximated by a finite difference of Gaussian functions: for a constant scale factor $k > 1$, he considers instead

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

Also, the practical implementation of [33] 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 [33]. For technical details we refer the reader to [54]. Here, and in the remaining 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

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

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 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 nor 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 confounded 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 Fig. 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 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 (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.

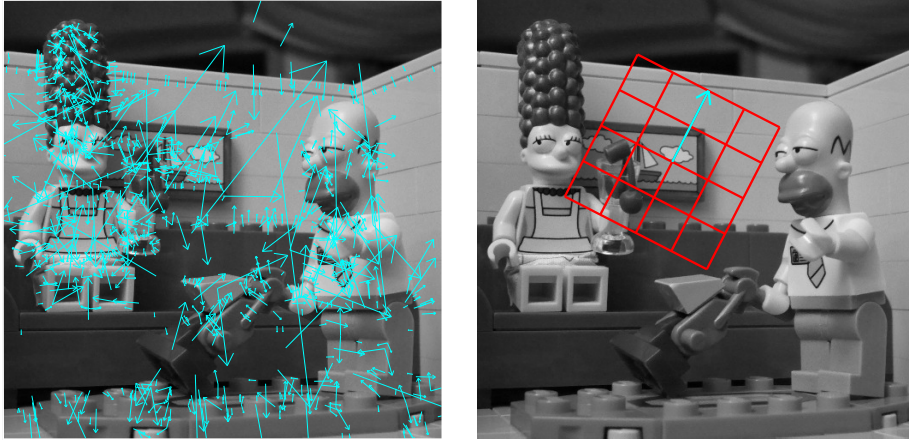


Figure 1: **Examples of SIFT keypoints and subcells.** On the left, one can see an original image (Courtesy of J. Delon) with overimposed SIFT oriented keypoints $(\mathbf{x}, \sigma, \alpha)$ represented as arrows originating from \mathbf{x} , with orientation α and length 6σ . On the right, we display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

3.1 Exponential Models with local HOG

We will denote by $\mathbb{T} = \mathbb{R}/2\pi\mathbb{Z}$ the set of angles, and \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 besides we adopt the same quantization bins for all SIFT subcells

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

(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

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

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

Maximum Entropy Distribution

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

$$\forall j \in \mathcal{J}, \forall \ell \in \{1, \dots, 8\}, \quad \mathbb{E}_P(f_{j,\ell}(\Theta)) = f_{j,\ell}(\theta_0), \quad (5)$$

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 in average the extracted statistics in the SIFT subcells, see Fig. 2.

There are many probability distributions P on \mathbb{T}^Ω that satisfy (5), and we will be mainly interested in the ones 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.

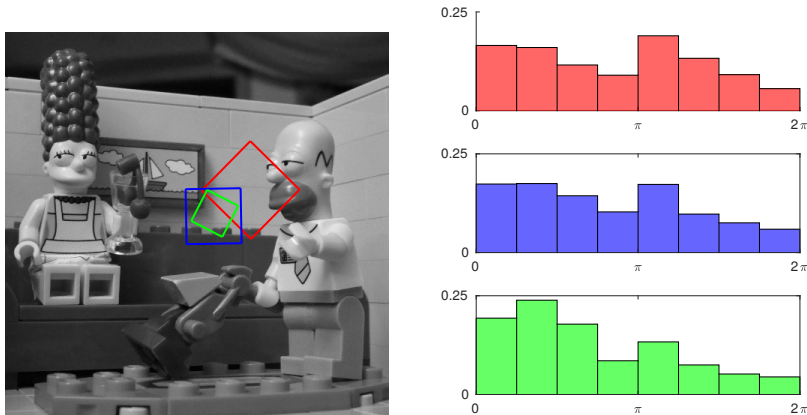


Figure 2: **Extracting HOG in SIFT subcells.** On the left, we display an original image (Courtesy of J. Delon) with three overlaid SIFT subcells s_j , and on the right, we display the corresponding HOG $(f_{j,\ell}(\theta_0))_{1 \leq \ell \leq 8}$ extracted in these subcells. The MaxEnt model is a probability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

Theorem 3.1. *There exists a family of numbers $\lambda = (\lambda_{j,\ell})_{j \in \mathcal{J}, 1 \leq \ell \leq 8}$ such that the probability distribution*

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

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 w.r.t. the Lebesgue measure $d\theta$ on \mathbb{T}^Ω satisfying the constraints (5).

Proof. This result directly follows from the general theorem given in [47]. The only difficulty is to handle 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). \square

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 $H(P)$:

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

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

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

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

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

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

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

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

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 to compute several statistics associated with this model. In particular, we can compute for any bounded measurable function $\psi : \mathbb{T} \rightarrow \mathbb{C}$

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

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

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

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

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

Notice that this system is highly non-linear and is in general difficult to solve.

A simple case: non-overlapped SIFT subcells

When a SIFT subcell s_j is not overlapped, then we have for any $\mathbf{x} \in s_j$, $|\mathcal{J}(\mathbf{x})| = 1$ and therefore

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

Then (14) gives

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

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

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

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 non-overlapped subcell s_j , there exists a constant $c_j > 0$ such that

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

Maximum-likelihood estimation

If the SIFT subcells intersect, there is no explicit solution anymore. 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

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

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

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

which gives

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

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

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

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 semi-positive 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 have 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 $L = 8|\mathcal{J}|$. Writing λ^k the iterates of the gradient descent with constant step size $h < \frac{2}{L}$, [49, Th 2.1.14] gives

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

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 to give more weight to the locations for which we have information at 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 2.1, 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 graylevels) 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

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

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 et al. [52] in order to copy pieces of an image into another one in a seamless way. More precisely, the goal is to minimize the functional

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

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

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

If we use periodic boundary conditions for the gradient, we can solve this problem with the Discrete Fourier Transform [44]. Indeed, if we use the simple derivation scheme based on periodic convolutions

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

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

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

Thus, for each $\boldsymbol{\xi}$ we have a barycenter problem which is simply solved by

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

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

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

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

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

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

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

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

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

3.3 Algorithm

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

For images having many SIFT keypoints in overlapping positions, this algorithm may be slow to converge as can be observed on the case of Fig. 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.

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 3.1 could be extended to statistics like

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

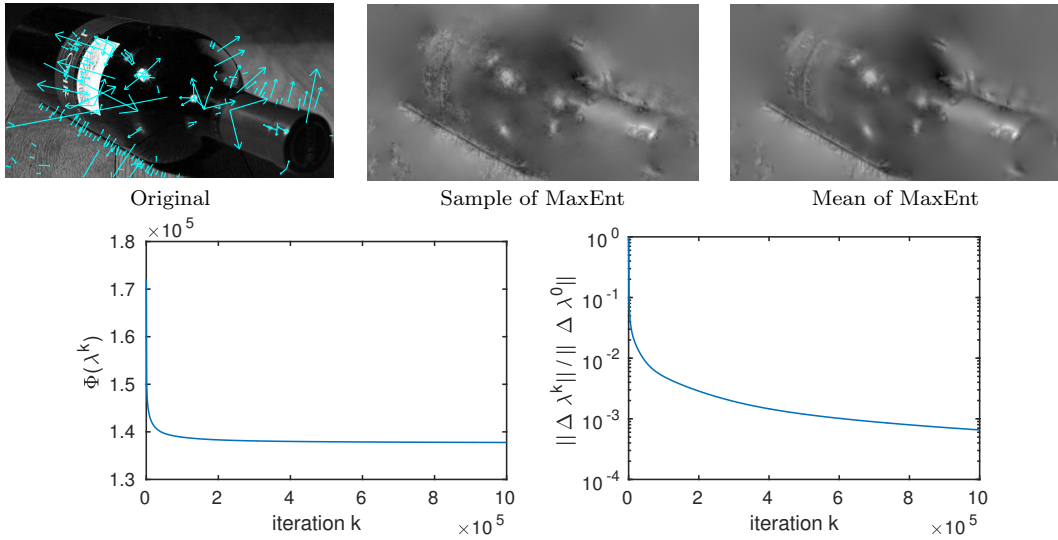


Figure 3: **Reconstruction with the MaxEnt model.** In the first row from left to right, we display an original image with overimposed 187 oriented keypoints, a sample of the associated MaxEnt model, and the expectation of the MaxEnt model. In the second row we display the evolution of Φ along the iterates, 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 3.2).

It is even possible to consider multiscale statistics using $\nabla g_{\sigma_j} * u$ instead of ∇u (as it will be the case in Section 4). But the analog of Proposition 3.2 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 untractable, and also only of minor interest. Indeed, in view of the results of Fig. 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

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

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

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

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 to set

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

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 paragraph.

4.2 Multiscale Poisson Reconstruction

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

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

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

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

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 $\mu = 50$.

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

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

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

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

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

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

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

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

As we will see in Section 5.1, the parameter μ allows 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 paragraph by summarizing the MS-Poisson sampling algorithm.

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

$$\forall \boldsymbol{\xi} \neq 0, \quad \widehat{\nu}_j(\boldsymbol{\xi}) = \frac{w(\sigma_j) \widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \widehat{\nabla}(\boldsymbol{\xi})^*}{\left(\mu + \sum_{j' \in \mathcal{J}} w(\sigma_{j'}) |\widehat{g}_{\sigma_{j'}}(\boldsymbol{\xi})|^2 \right) |\widehat{\nabla}(\boldsymbol{\xi})|^2} \quad \text{and } \widehat{\nu}_j(0) = 0. \quad (43)$$

Then, as in Section 3.2 we get the convolutive expression

$$U = \sum_{j \in \mathcal{J}} \nu_j * V_j = \sum_{j \in \mathcal{J}} \left(\nu_{j,1} * V_{j,1} + \nu_{j,2} * V_{j,2} \right). \quad (44)$$

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

$$\mathbb{E}(U) = \sum_{j \in \mathcal{J}} \nu_j * \mathbb{E}(V_j), \quad (45)$$

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

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

$$\text{Var}(U(\mathbf{x})) = \sum_{j \in \mathcal{J}} \text{Var}(\nu_j * V_j(\mathbf{x})). \quad (46)$$

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

$$\text{Var}(\nu_j * V_j(\mathbf{x})) = \text{Var} \left(\sum_{\mathbf{y} \in \Omega} \nu_j(\mathbf{x} - \mathbf{y}) V_j(\mathbf{y}) \right) = \sum_{\mathbf{y} \in \Omega} \text{Var}(\nu_j(\mathbf{x} - \mathbf{y}) V_j(\mathbf{y})) \quad (47)$$

$$= \sum_{\mathbf{y} \in \Omega} \nu_j(\mathbf{x} - \mathbf{y}) \text{Cov}(V_j(\mathbf{y})) \nu_j^T(\mathbf{x} - \mathbf{y}) \quad (48)$$

$$= \sum_{\mathbf{y} \in \Omega} \nu_{j,1}^2(\mathbf{x} - \mathbf{y}) \text{Var}(V_{j,1}(\mathbf{y})) + \nu_{j,2}^2(\mathbf{x} - \mathbf{y}) \text{Var}(V_{j,2}(\mathbf{y})) \quad (49)$$

$$+ 2\nu_{j,1}(\mathbf{x} - \mathbf{y}) \nu_{j,2}(\mathbf{x} - \mathbf{y}) \text{Cov}(V_{j,1}(\mathbf{y}), V_{j,2}(\mathbf{y})). \quad (50)$$

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

$$\text{Cov}(U(\mathbf{x}), U(\mathbf{y})) = \sum_{j \in \mathcal{J}} \sum_{\mathbf{z} \in \Omega} \nu_j(\mathbf{x} - \mathbf{z}) \text{Cov}(V_j(\mathbf{z})) \nu_j^T(\mathbf{y} - \mathbf{z}). \quad (51)$$

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. 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

Let us first compare the reconstruction results obtained with MaxEnt and with MS-Poisson. On Fig. 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 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 on Fig. 6.

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 the ones obtained with MaxEnt. Besides, 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 term is also a source of blur in the result. This is confirmed by Fig. 5 where we display several MS-Poisson reconstructions with varying regularization parameter μ . In Fig. 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 Fig. 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 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 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 its t-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 e.g.), 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 Fig. 4, it is obvious if one tries to set the gradient magnitude to 1 in the global Poisson reconstruction.

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 Fig. 6 we display the standard deviations of all pixels values in each model. One first remark is that MaxEnt has in general 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

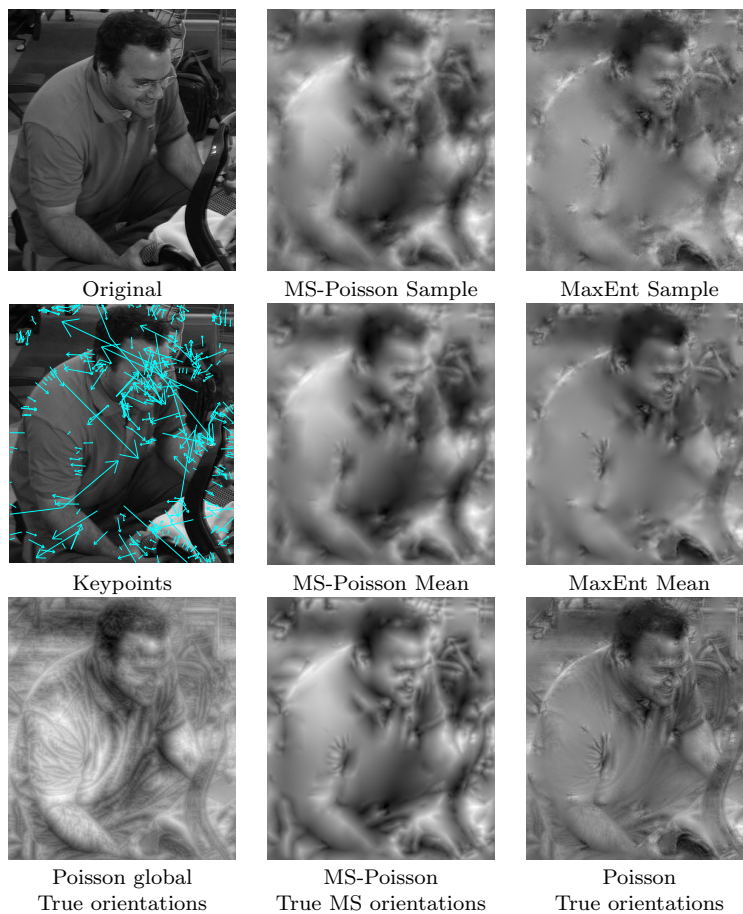


Figure 4: **Reconstruction results with MaxEnt and MS-Poisson models.** In the first column we display an original image, the corresponding oriented keypoints, and the Poisson reconstruction with true gradient orientations of the whole image and magnitude set to 1. In the second column we display a sample 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. See the text for comments on these results. (Images are better seen on the electronic version)



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. We empirically observed that a good compromise between recovered details and smoothness is often attained around $\mu = 50$. (Images are better seen on the electronic version)

details. That being said, the variance of both these models is relatively small compared to the global range of the mean image, which indicates that both these models have quite small variations around the mean.

Let us emphasize that in our experiments, we used all the keypoints computed by the SIFT methods and we 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 still, it is clear that for some images, the reconstruction will be quite different when discarding those keypoints. For example in the case of Fig. 7, 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.

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 on Fig. 8, we get back similar keypoints in many regions of the image, but still with some variations in positions, scales and orientations. 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 less keypoints in the MS-Poisson model: indeed, since it is more regular we lose some extrema in the scale-space. Besides, the regularization tends to change the scale of the structures, 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 [54] (we used the proposed default parameters). This algorithm follows the matching method proposed in [33] which essentially pairs SIFT keypoints by thresholding the ratio between the distances to the first and second nearest neighbors (computed with the ℓ^2 -

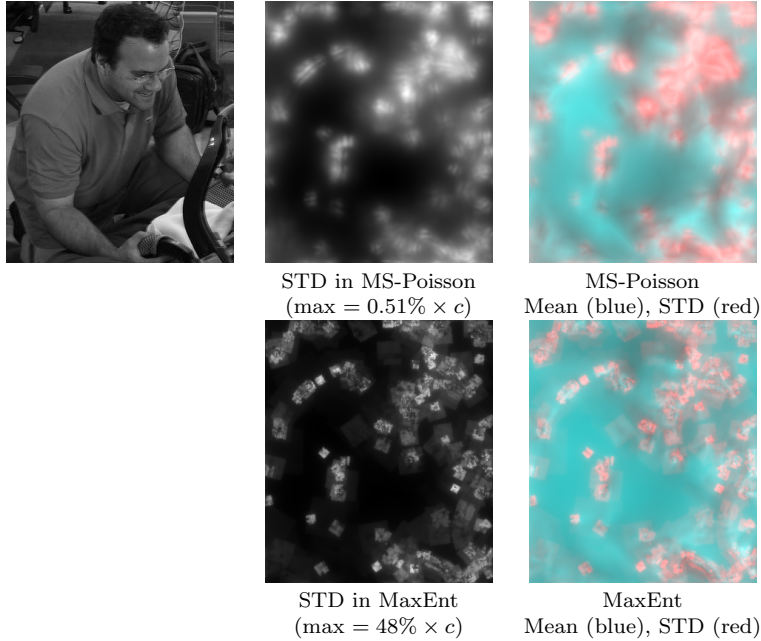


Figure 6: **Standard deviations of MS-Poisson and MaxEnt models.** On the top left we display the original image. On the rest of the figure we display the images formed with the standard deviations (STD) of the models MS-Poisson (first row) and MaxEnt (second row). On the second column we display the raw STD values. On 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 on the electronic version)

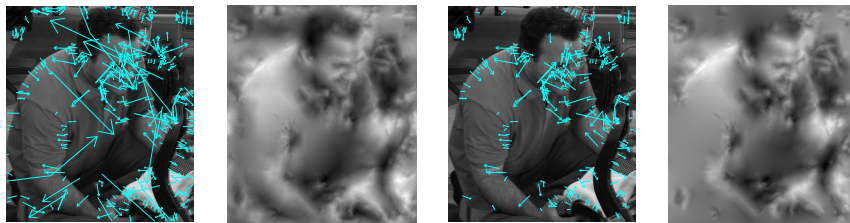


Figure 7: **Discard keypoints near image boundary.** In this figure, we examine the effect of discarding keypoints whose associated SIFT cell is not entirely contained in the image domain. The displayed reconstructions are samples of the MS-Poisson model.

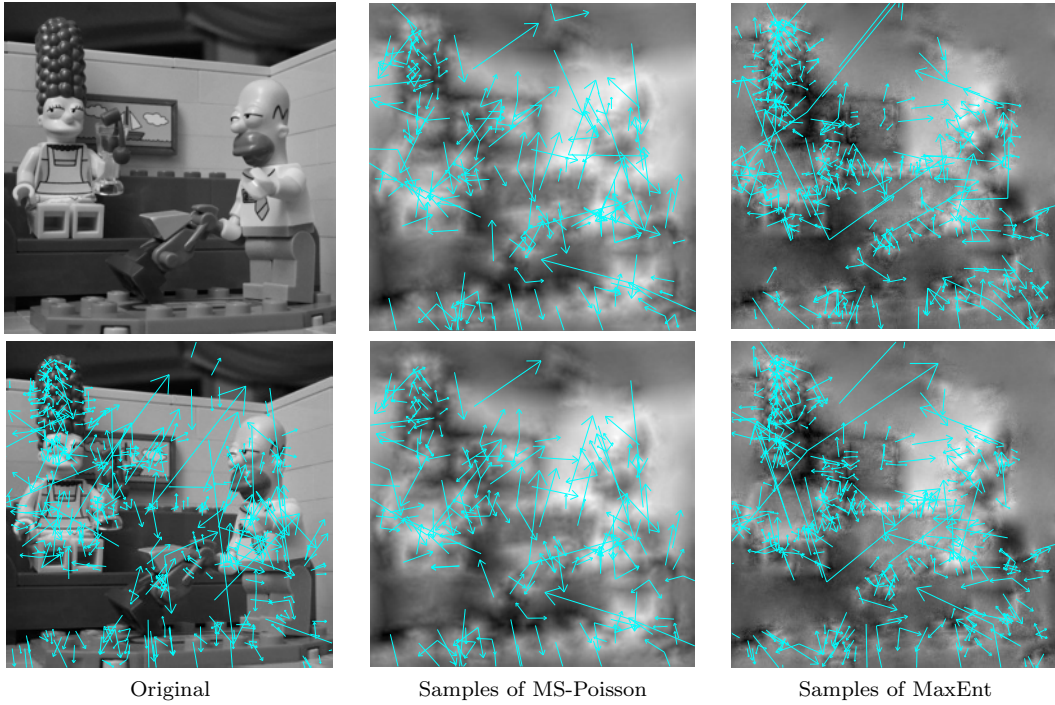


Figure 8: **Keypoints after reconstruction.** In the first column we display an original image and the same image with its SIFT keypoints. In the second column we display two samples of the MS-Poisson model. In the third column we display two samples of the MaxEnt model. We display the keypoints associated to these images as overimposed blue arrows. Notice that several keypoints are retrieved after reconstruction, with still some variations in positions and orientations. 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 on the electronic version)

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 Fig. 8, 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 less keypoints are correctly matched: over 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 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

$$\tilde{H}_{j,\ell} = \frac{f_{16(p-1)+\ell}}{\sum_{\ell'=1}^8 f_{16(p-1)+\ell'}} \quad (1 \leq \ell \leq 8). \quad (52)$$

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

On Fig. 9, we display several reconstruction results obtained with the model MS-Poisson 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 the ones 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 Equation (2)), which dramatically reduce the quantity of information. Also, in the true SIFT descriptors, the pixels vote for orientations 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 paragraph, we propose to compare our reconstruction models with the ones obtained by the methods by Weinzaepfel et al. [64] and Dosovitskiy & Brox [19]. 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 [36] 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).

They are also minor differences in the extracted information because both 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 [42]) 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 in 4×4 cells and each cell containing a keypoint is being associated with the corresponding oriented keypoint

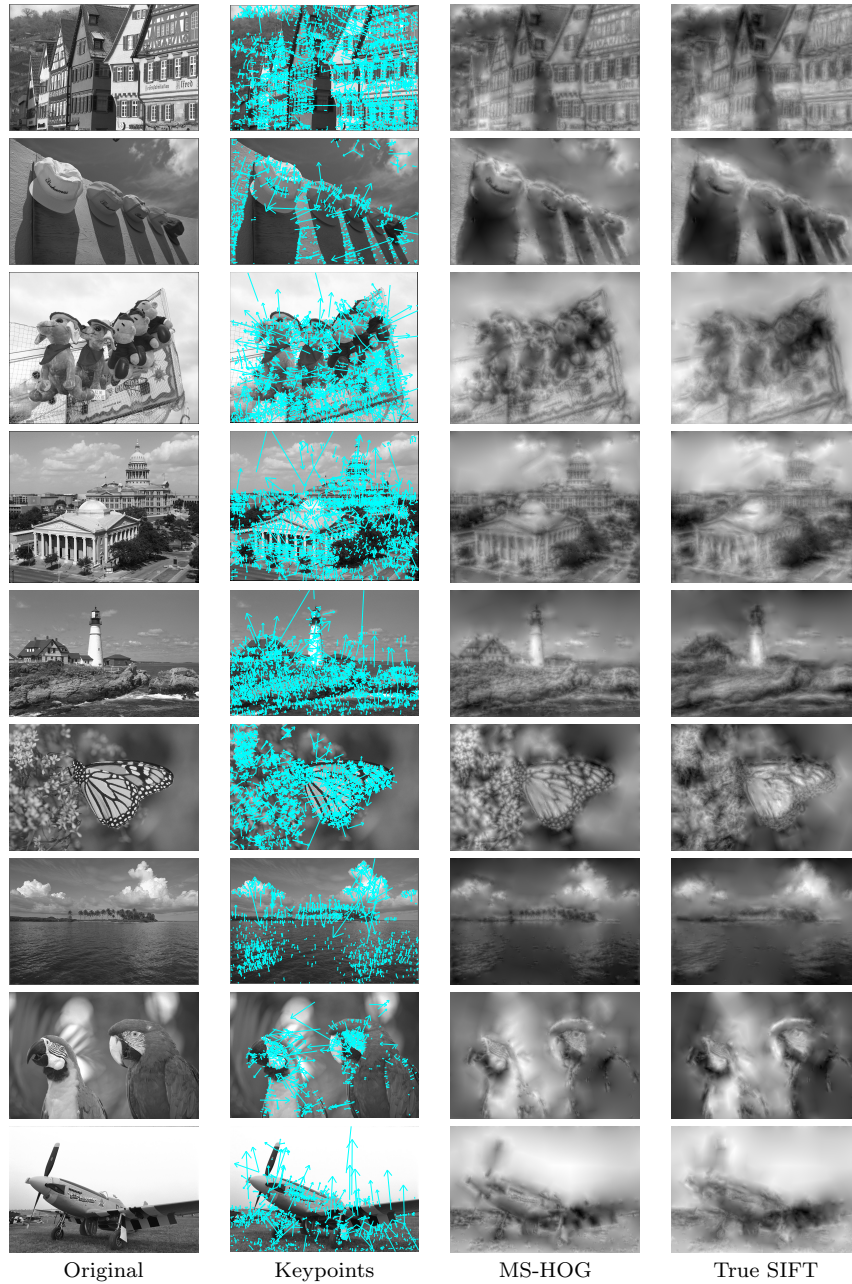


Figure 9: **Reconstruction results from multiscale HOG or SIFT descriptors** with images of the Live database [57]. For each row, from left to right, we display an original image, the same image with overimposed SIFT keypoints, a sample of the MS-Poisson model obtained from multiscale HOG, and a sample 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.

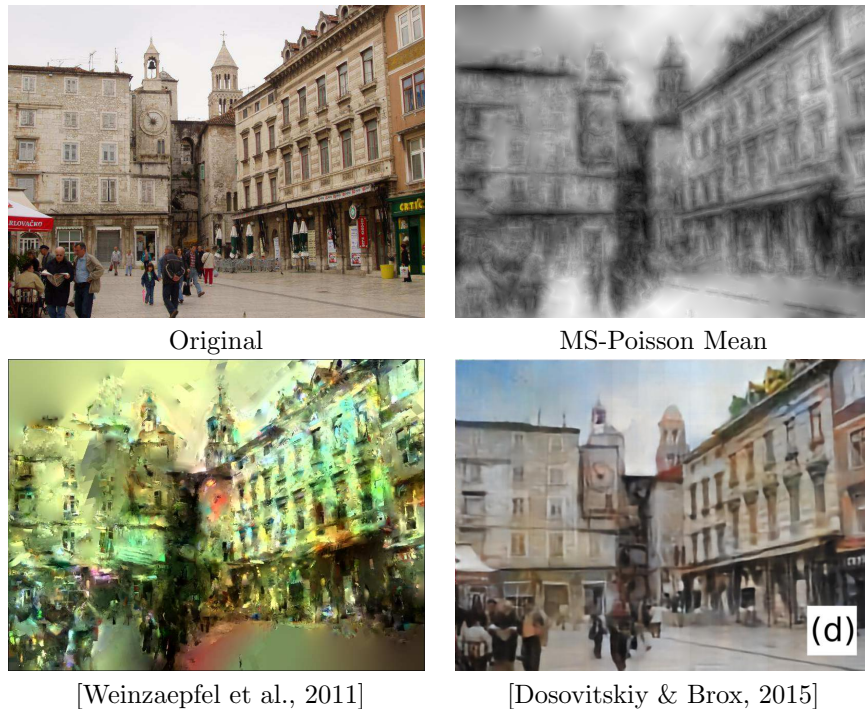


Figure 10: **Comparison for SIFT reconstruction.** In the first row we display the original image and the reconstruction results obtained as the expectation of the MS-Poisson model computed on the true SIFT descriptors (see Section 5.2). In the second row we display the results obtained with the methods 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.

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 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 the one of [19], and in particular it cannot get back information which are completely lost in the SIFT descriptors (global contrast, or also color information).

5.4 Reconstruction with other keypoints

In this paragraph we question the very definition of the SIFT keypoints in terms of synthesis, in a similar way that what was done in [50]. 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 to compare 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 Taylor formula around a point \mathbf{x} , one can write when $\sigma \rightarrow 0$ that

$$\int u(\mathbf{x} + \mathbf{z})g_\sigma(\mathbf{z})d\mathbf{z} - u(\mathbf{x}) = \sigma^2 \Delta u(\mathbf{x}) + o(\sigma^2). \quad (53)$$

Therefore, nearby 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, and 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

$$(\mathbf{x}, \sigma) \mapsto |g_\sigma * u(\mathbf{x}) - u(\mathbf{x})|. \quad (54)$$

In our implementation, we detect these maxima on a discretized scale-space with 30 scales $s = 2^{r/6}$, $0 \leq r < 30$. Besides, in order to draw a comparison with a fixed number of keypoints, we only keep the 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 than the ones provided by the SIFT method.

The second and third choices (“Random-unif” and “Random-grad”) consists 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. Modelling 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 on Fig. 11. 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 Fig. 11 the average scale of these keypoints is approximately twice larger than the one of the SIFT keypoints. Besides, 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 the one 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, both from the analysis or the synthesis perspective.

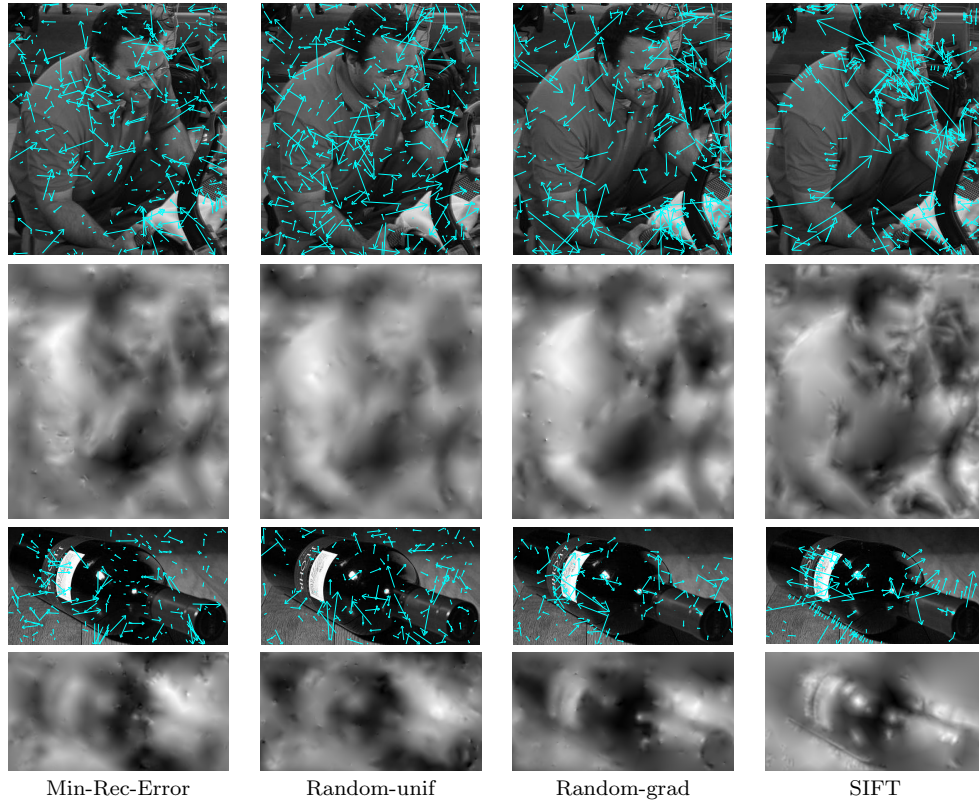


Figure 11: **Reconstruction with other keypoints.** The first column (“Minimum reconstruction error”) corresponds to the keypoints obtained as local minima of (54). The second (“Random-unif”) and third column (“Random-grad”) corresponds to the randomly selected keypoints. The last column corresponds to the standard SIFT keypoints. The original images are displayed on Fig. 3 and Fig. 4. See the text in Section 5.4 for the precise definition of these sets of keypoints, and additional comments.

6 Conclusion

In this paper we proposed two stochastic models (MaxEnt, respectively MS-Poisson) for reconstructing an image based only on the information contained in the (monoscale, respectively 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 models over competing approaches is that they do not rely on any external image database, and besides the convolutive expressions found in this paper allow to compute 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 transactions on pattern analysis and machine intelligence, 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*, arXiv preprint arXiv:1401.0542, (2015).
- [4] B. ALLEN AND M. KON, *Unique recovery from edge information*, in Sampling Theory and Applications (SampTA), 2015 International Conference on, IEEE, 2015, pp. 312–316.
- [5] F. ATTNEAVE, *Some informational aspects of visual perception.*, Psychological review, 61 (1954), p. 183.
- [6] S. BATTIATO, G. GALLO, G. PUGLISI, AND S. SCCELLATO, *SIFT features tracking for video stabilization*, in Image Analysis and Processing, 2007. ICIAP 2007. 14th International Conference on, IEEE, 2007, pp. 825–830.
- [7] H. BAY, A. ESS, T. TUYTELAARS, AND L. VAN GOOL, *Speeded-up robust features (SURF)*, Computer vision and image understanding, 110 (2008), pp. 346–359.
- [8] M. BLACK AND A. JEPSON, *Eigenttracking: Robust matching and tracking of articulated objects using a view-based representation*, International Journal of Computer Vision, 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 Computer Vision (ICCV), 2011 IEEE International Conference on, 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. Am. A, 4 (1987), pp. 221–231.
- [13] S. CURTIS, S. SHITZ, AND A. OPPENHEIM, *Reconstruction of nonperiodic two-dimensional signals from zero crossings*, IEEE Transactions on Acoustics, Speech, and Signal Processing, 35 (1987), pp. 890–893.
- [14] N. DALAL AND B. TRIGGS, *Histograms of oriented gradients for human detection*, in Proceedings of the IEEE CVPR, vol. 1, 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 Transactions on PAMI, 36 (2014), pp. 874–887.
- [16] A. DESOLNEUX, *When the a contrario approach becomes generative*, International Journal of Computer Vision, 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), Springer, Lecture Notes in Computer Science, 2017, pp. 133–145.
- [18] A. DESOLNEUX, L. MOISAN, AND J. MOREL, *From Gestalt theory to image analysis: a probabilistic approach*, vol. 34, Springer Science & Business Media, 2007.
- [19] A. DOSOVITSKIY AND T. BROX, *Inverting Visual Representations with Convolutional Networks*, arXiv:1506.02753 [cs], (2015).
- [20] J. H. ELDER AND S. W. ZUCKER, *Scale space localization, blur, and contour-based image coding*, in Computer Vision and Pattern Recognition, 1996. Proceedings CVPR’96, 1996 IEEE Computer Society Conference on, 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 Transactions on PAMI, 32 (2010), pp. 1627–1645.
- [23] C. HARRIS AND M. STEPHENS, *A combined corner and edge detector.*, in Alvey vision conference, vol. 15, Citeseer, 1988, p. 50.
- [24] R. HUMMEL AND R. MONIOT, *Reconstructions from zero crossings in scale space*, IEEE Transactions on Acoustics, Speech, and Signal Processing, 37 (1989), pp. 2111–2130.
- [25] F. JUEFEI-XU AND M. SAVVIDES, *Learning to invert local binary patterns.*, in BMVC, 2016.

- [26] S. K. AND Z. A., *Very deep convolutional networks for large-scale image recognition*, in Proceedings of the International Conference on Learning Representations, 2014.
- [27] H. KATO AND T. HARADA, *Image reconstruction from bag-of-visual-words*, in Proceedings of the IEEE CVPR, 2014, pp. 955–962.
- [28] 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.
- [29] 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.
- [30] 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.
- [31] T. LINDBERG, *Feature detection with automatic scale selection*, International journal of computer vision, 30 (1998), pp. 79–116.
- [32] T. LINDBERG, *Image matching using generalized scale-space interest points*, Journal of Mathematical Imaging and Vision, 52 (2015), pp. 3–36.
- [33] D. LOWE, *Distinctive image features from scale-invariant keypoints*, International Journal of Computer Vision, 60 (2004), pp. 91–110.
- [34] Y. LU, S. ZHU, AND Y. N. WU, *Learning FRAME models using CNN filters for knowledge visualization*, CoRR, abs/1509.08379 (2015), <http://arxiv.org/abs/1509.08379>, arXiv:1509.08379.
- [35] A. MAHENDRAN AND A. VEDALDI, *Understanding deep image representations by inverting them*, in IEEE CVPR, 2015, pp. 5188–5196.
- [36] A. MAHENDRAN AND A. VEDALDI, *Visualizing deep convolutional neural networks using natural pre-images*, International Journal of Computer Vision, 120 (2016), pp. 233–255.
- [37] 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.
- [38] S. MALLAT AND S. ZHONG, *Characterization of signals from multiscale edges*, IEEE Transactions on PAMI, 14 (1992), pp. 710–732.
- [39] D. MARR, *Vision: A computational investigation into the human representation and processing of visual information*, W.H. Freeman and Company, 1982.
- [40] D. MARR AND E. HILDRETH, *Theory of edge detection*, Proceedings of the Royal Society of London B: Biological Sciences, 207 (1980), pp. 187–217.
- [41] Y. MEYER, *Wavelets-algorithms and applications*, vol. 1, Society for Industrial and Applied Mathematics Translation, 1993.
- [42] K. MIKOLAJCZYK AND C. SCHMID, *Scale & affine invariant interest point detectors*, International journal of computer vision, 60 (2004), pp. 63–86.
- [43] K. MIKOLAJCZYK AND C. SCHMID, *A performance evaluation of local descriptors*, IEEE Transactions on PAMI, 27 (2005), pp. 1615–1630.

- [44] J.-M. MOREL, A. PETRO, AND C. SBERT, *Fourier implementation of Poisson image editing*, Pattern Recognition Letters, 33 (2012), pp. 342–348.
- [45] J.-M. MOREL AND G. YU, *ASIFT: A new framework for fully affine invariant image comparison*, SIAM Journal on Imaging Sciences, 2 (2009), pp. 438–469.
- [46] J.-M. MOREL AND G. YU, *Is SIFT scale invariant?*, Inverse Problems and Imaging, 5 (2011), pp. 115–136.
- [47] D. MUMFORD AND A. DESOLNEUX, *Pattern Theory: The Stochastic Analysis of Real-World Signals*, A K Peters/CRC Press, Natick, Mass, 2010.
- [48] P. MUSÉ, F. SUR, F. CAO, Y. GOUSSEAU, AND J.-M. MOREL, *An a contrario decision method for shape element recognition*, International Journal of Computer Vision, 69 (2006), pp. 295–315.
- [49] Y. NESTEROV, *Introductory lectures on convex optimization: A basic course*, vol. 87, Springer, 2004.
- [50] M. NIELSEN AND M. LILLHOLM, *What do features tell about images?*, in Scale-Space, vol. 1, Springer, 2001, pp. 39–50.
- [51] T. OJALA, M. PIETIKÄINEN, AND T. MÄENPÄÄ, *Multiresolution gray-scale and rotation invariant texture classification with local binary patterns*, IEEE Transactions on PAMI, 24 (2002), pp. 971–987.
- [52] P. PÉREZ, M. GANGNET, AND A. BLAKE, *Poisson Image Editing*, in ACM SIGGRAPH 2003 Papers, SIGGRAPH '03, 2003, pp. 313–318, doi:10.1145/1201775.882269.
- [53] J. PHILBIN, O. CHUM, M. ISARD, J. SIVIC, AND A. ZISSERMAN, *Object retrieval with large vocabularies and fast spatial matching*, in IEEE Conference on Computer Vision and Pattern Recognition, 2007., IEEE, 2007, pp. 1–8.
- [54] I. REY OTERO AND M. DELBRACIO, *Anatomy of the SIFT Method*, Image Processing On Line, 4 (2014), pp. 370–396, doi:10.5201/ipo1.2014.82.
- [55] E. ROSTEN AND T. DRUMMOND, *Machine learning for high-speed corner detection*, in Computer Vision–ECCV 2006, Springer, 2006, pp. 430–443.
- [56] J. SANZ AND T. HUANG, *Theorems and experiments on image reconstruction from zero crossings*, 1987. IBM Almaden Research Center.
- [57] H. SHEIKH, Z. WANG, L. CORMACK, AND A. BOVIK, *Live image quality assessment database release 2 (2005)*, 2005.
- [58] J. SIVIC AND A. ZISSERMAN, *Video Google: A text retrieval approach to object matching in videos*, in Proceedings of the IEEE ICCV, 2003, pp. 1470–1477.
- [59] T. TUYTELAARS AND K. MIKOLAJCZYK, *Local invariant feature detectors: a survey*, Foundations and trends in computer graphics and vision, 3 (2008), pp. 177–280.
- [60] T. TUYTELAARS AND L. VAN GOOL, *Matching widely separated views based on affine invariant regions*, International journal of computer vision, 59 (2004), pp. 61–85.
- [61] A. VEDALDI AND B. FULKERSON, *Vfeat: 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 IEEE ICCV, 2013, pp. 1–8.
- [63] C. WALLRAVEN, B. CAPUTO, AND A. GRAF, *Recognition with local features: the kernel recipe*, in Proceedings of the IEEE ICCV, 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 Transactions on Circuits and Systems for Video Technology, 19 (2009), pp. 945–954.
- [66] A. YILMAZ, O. JAVED, AND M. SHAH, *Object tracking: A survey*, ACM computing surveys (CSUR), 38 (2006), p. 13.
- [67] M. ZEILER AND R. FERGUS, *Visualizing and understanding convolutional networks*, in Proceedings of ECCV, 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*, International journal of computer vision, 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*, International Journal of Computer Vision, 27 (1998), pp. 107–126.