

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 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.

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,

1. Introduction. ¹

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 transformations² (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 to [43, 59, 45, 32] for a more comprehensive survey.

Harris and Stephens proposed a combined corner and edge detector based on the

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

¹A preliminary version of this work was published as a conference paper in [17].

²The translation invariance is generally always required, and often trivial.

45 determinant and trace of the structure tensor of the image [23]. A multiscale variant
 46 based on a normalized Laplacian of Gaussian (LoG) scale-space, coined Harris-Laplace
 47 was proposed by Mikolajczyk and Schmid [42]. The same authors also proposed in [42]
 48 the Harris-affine point detector which extends the previous one with a normalization
 49 step in order to get invariance to affine transformations. Tuytelaars and Mikolajczyk
 50 proposed in [60] two region detectors both starting from anchor points (e.g. Harris
 51 points); then the first one selects a region within detected edges around the anchor,
 52 and the second one extracts a region by analyzing intensity profiles on rays emanating
 53 from the anchor. Rosten and Drummond introduced in [55] the “features from ac-
 54 celerated segment test” (FAST) which is a corner detector accelerated by a machine
 55 learning technique. This approach has been further fastened by Mair et al. [37] using
 56 optimal decision trees, thus obtaining an “adaptive and generic accelerated segment
 57 test” (AGAST). Musé et al. proposed in [48] to extract shapes from the image level
 58 lines, and to process them in order to get an affine invariant representation.

59 In parallel of this research on keypoints, many techniques have been proposed
 60 for invariant local descriptions of images. An early descriptor is given by the local
 61 binary patterns (LBP) defined by Ojala et al. [51] which extracts signs of differences
 62 of image values on pixels located on a circular neighborhood of a keypoint. The LBP
 63 were originally designed for texture description but can also be used for face detec-
 64 tion [1]. In [33], Lowe introduced the scale invariant feature transform (SIFT) which
 65 first extracts the keypoints as local extrema of the “Difference of Gaussian” (DoG)
 66 approximation of the LoG, and next computes around each keypoint a local descriptor
 67 based on normalized histograms of gradient direction (HOG), see the details in Sec-
 68 tion 2. Notice that similar HOG descriptors computed on a dense grid were actually
 69 used in [14] for person detection; one reference implementation of the HOG descrip-
 70 tors is given in [22]. A fully affine-invariant extension of SIFT, named ASIFT, was
 71 proposed by Morel and Yu [45] and consists in applying the SIFT method with the
 72 image transformed with several simulated affine maps. The SURF method (Speeded-
 73 up robust features) proposed by Bay et al. [7] is closely related in construction to the
 74 SIFT method, but allows for a faster implementation. At a higher semantic level,
 75 local image behavior can be also represented as visual words [58, 11] which are ob-
 76 tained as cluster points in a feature space. Later, some authors proposed to describe
 77 a patch using local binary descriptors (LBD), which extracts the signs of differences
 78 between Gaussian measurements taken at different locations. Using different ways of
 79 selecting these locations leads to the methods BRISK [30] (binary robust invariant
 80 scalable keypoints) or FREAK [2] (fast retina keypoint). All of these descriptors have
 81 quite different invariance properties (evaluated either in a theoretical or experimental
 82 framework).

83 Long before the design of these image descriptors, the question of a minimal
 84 representation of an image was thoroughly studied, mainly for compression purpose.
 85 Through the concept of *raw primal sketch*, Marr [39] suggested that the human visual
 86 system processes images by retaining essentially the lines of zero-crossing of the Lapla-
 87 cian at several scales. This leads to the conjecture that an image is uniquely defined
 88 by these zero-crossing lines, a conjecture that was later precised by Mallat [38] using
 89 wavelet modulus maxima. Both these conjectures were proved wrong by Meyer [41]
 90 but still, algorithms for approximate reconstruction were proposed by Hummel and
 91 Moniot [24] for zero-crossings and by Mallat and Zhong [38] for the case of wavelet
 92 modulus maxima. Besides, unique characterization can be shown to be true under
 93 some additional hypotheses [12, 13, 56, 4, 3].

94 From a more practical point of view, several authors have raised the question of

95 inversion of a feature-based representation. For example, Elder and Zucker [20] pro-
 96 posed an algorithm for image reconstruction from detected contours, based on the heat
 97 diffusion. Nielsen and Lillholm [50] consider the problem of variational reconstruc-
 98 tion from linear measurements; in addition to the minimum variance reconstruction
 99 (given by the pseudo-inverse of the measurements matrix), they propose two varia-
 100 tional reconstructions based on either the entropy (of the image seen as a probability
 101 distribution on its domain) or the H^1 norm. Interestingly, they discuss the problem
 102 of extracting a subset of linear measurements which leads to the best reconstruction
 103 and empirically compare three different strategies for that purpose.

104 Motivated by privacy issues (since the descriptors may be transmitted on an
 105 unsecured network), Weinzaepfel et al. [64] addressed image reconstruction from the
 106 output of a SIFT transform adapted with elliptic keypoints. One important difference
 107 with previous works is that this method exploits a database of image patches: for
 108 each keypoint, a patch with similar description is looked for in the database, and
 109 all the patches are stitched together with Poisson image editing [52]. Vondrick et
 110 al. [62] address reconstruction from dense HOGs by relying on a paired dictionary
 111 representation of HOGs and patches. Also, d’Angelo et al. [15] address reconstruction
 112 from local binary descriptors by relying on primal-dual optimization techniques; in
 113 contrast with [64, 62], this method does not need any external information. Kato
 114 and Harada [27] formulate reconstruction from bag of visual words as a problem of
 115 quadratic assignment. Finally, Juefei-Xu and Savvides [25] propose to invert the
 116 LBP representation with an approach based on paired dictionary learning with an ℓ^0
 117 constraint.

118 More recently, the success of deep convolutional neural networks in image classi-
 119 fication [28, 67] has urged the need of inverting the corresponding representations in
 120 order to intuitively understand the kind of information that is extracted at each layer.
 121 Even if they do not formulate it as an inverting procedure, Zeiler and Fergus [67] pro-
 122 posed to build a deconvolution network that allows to visualize in image space the
 123 stimuli that excite one response at a particular layer of the neural network. Given an
 124 image u , Mahendran and Vedaldi [35, 36] proposed to search for a pre-image of an
 125 image representation $\varphi(u)$ by minimizing a functional containing a loss term related
 126 to the representation φ and a regularizing term (in particular the H^1 norm). Even if
 127 the regularizer is convex, the transformation φ is in general highly non-linear so that
 128 the resulting optimization problem is not convex; so the output of the inversion may
 129 depend on the parameters and initializations of the chosen optimization procedure.
 130 On the other hand, Dosovitskiy and Brox [19] suggest to learn an approximate left in-
 131 verse of the representation (i.e. a mapping φ_L^{-1} such that $\varphi_L^{-1}(\varphi(u)) \approx u$ for every u)
 132 in the form of an up-convolutional network. These methods are generic in the sense
 133 that they can be applied to any image representation that can be approximated by
 134 the output of a convolutional neural network; in particular, the authors of [19] display
 135 inversion results for both HOG, SIFT and AlexNet [28] representations. Notice that
 136 the inversion/visualization techniques of [67, 19] exploit an external database while
 137 the one of [35, 36] does not.

138 Instead of building a uniquely defined inversion technique (using regularization),
 139 another way to perform reconstruction from the image representation φ is to sample
 140 from a stochastic model that explores the set of pre-images of $\varphi(u)$. This is particu-
 141 larly relevant if one uses an image representation that is not invertible: for example,
 142 the SIFT cells of an image may not cover its whole domain and thus many images
 143 could have the same SIFT descriptors. Besides, the HOG descriptors are inherently
 144 of a statistical nature: each HOG extracts the distribution of gradient orientations

145 in one small area. Thus they only provide a locally pooled information and thus do
 146 not precisely constrain each gradient value. For this reason, the inversion by direct
 147 (regularized) optimization proposed in [35, 36] is not adapted to the usual SIFT rep-
 148 resentation (sometimes called sparse SIFT as opposed to SIFT descriptors computed
 149 on a dense grid).

150 One way to address this problem is to sample from the maximum entropy model
 151 that complies with these statistical constraints. Such maximal entropy models were
 152 considered by Zhu, Wu and Mumford in [69, 47] for texture modelling based on re-
 153 sponses to an automatically selected subset of filters chosen in a filter bank. This
 154 approach has been recently extended by Lu, Zhu and Wu to responses to a pre-
 155 trained neural network [34]. Maximum entropy models were also used to question the
 156 noise models used in the *a contrario* framework for feature detections in images [18]:
 157 in [16], for two types of given detections (cluster of points, or line segments), Desol-
 158 neux proposes explicit computations of maximal entropy image models that lead to
 159 the same detections (in average). Let us emphasize that one important difference
 160 with previous works is that, more than reconstructing the original image, we aim at
 161 exploring the set of images with similar HOG description at the keypoints positions,
 162 with the least possible *a priori* of what the reconstruction should look like. In con-
 163 trast, the dependence on an external database in [52, 19] poses a strong *a priori* on
 164 the reconstruction.

165 In the present paper, we propose two stochastic models that complies with sta-
 166 tistical features given by a SIFT-like representation. In order to derive explicit com-
 167 putations, we work on a simplified SIFT transform which extracts multiscale HOGs
 168 from regions around the (usual) SIFT keypoints. The first model, called MaxEnt, is
 169 indeed an instance of maximum entropy model which complies with local statistical
 170 constraints on the gradient orientations (at scale 0, i.e. the image scale). Once the
 171 parameters of this model are estimated (using a gradient descent), a target gradient
 172 orientation can be sampled, and we recover an image by solving a classical Poisson
 173 problem. The second model, called MS-Poisson, consists in first independent sam-
 174 pling of multiscale gradient orientations in all the SIFT cells, and next merging all
 175 the pieces by solving a global multiscale Poisson problem. Even if this model does not
 176 solve an explicit maximum entropy problem, it allows to coherently merge information
 177 given at several scales. Several experiments show that both these models are able to
 178 recover large image structures and compare well to the results of [64] while not using
 179 any external information. Finally, we discuss the definition of the SIFT keypoints in
 180 terms of optimality of reconstruction, thus raising the following question related to
 181 visual information theory: “Can we measure the optimality (at fixed memory budget)
 182 of some image descriptor in terms of reconstruction?”

183 The paper is organized as follows. In Section 2, we briefly recall the main steps
 184 of the SIFT method, and explain the simplified SIFT descriptors that we use for
 185 reconstruction. In Section 3, we build and study the maximum entropy model (Max-
 186 Ent) used for reconstruction from monoscale HOGs computed in the SIFT subcells.
 187 In Section 4, we propose the multiscale Poisson model (MS-Poisson) that allows to
 188 comply with multiscale HOGs taken in the SIFT subcells; the corresponding H^1 -
 189 regularized multiscale Poisson problem is explicitly solved. Finally, in Section 5 we
 190 display several reconstruction results obtained with both models (applied with sim-
 191 plified SIFT, or also the true SIFT), study the variability of the reconstruction (in
 192 terms of first and second order moments, but also of SIFT keypoints computed on the
 193 reconstruction). We also compare with other existing reconstruction techniques and
 194 apply the reconstruction models on other keypoint sets, thus confirming (from the

195 synthesis perspective) the efficiency of the SIFT method for global image description.
 196 Finally in Section 6 we conclude the discussion proposed in this paper and open some
 197 perspectives for future research. A preliminary version of this work was published as
 198 a conference paper in [17]. Compared to the conference version, here we explain in
 199 more details the derivation of MaxEnt and MS-Poisson models and we provide some
 200 more properties of these models and in particular explicit formulae for the first and
 201 second order moments of these models. We also propose several new experiments
 202 which illustrate the performance and the variability of these models (with qualitative
 203 and quantative evaluation) and question the role of the keypoint definition in the
 204 quality of reconstruction.

205 **2. A Brief Summary of the SIFT Method.** In this section we briefly recall
 206 the construction of keypoints and local descriptors used in the SIFT method, and we
 207 explain the simplified descriptors that will be later used for the reconstruction in the
 208 next sections.

209 **2.1. Gaussian Scale-Space and Keypoints.** Following [31], we introduce the
 210 Gaussian scale-space in a continuous domain. Let $u : \mathbb{R}^2 \rightarrow \mathbb{R}$ be an integrable
 211 function. For $\sigma > 0$, we introduce the function $g_\sigma : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by

$$212 \quad g_\sigma(\mathbf{x}) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{|\mathbf{x}|^2}{2\sigma^2}\right).$$

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

$$214 \quad \forall \mathbf{x} \in \mathbb{R}^2, \forall \sigma > 0, \quad L_u(\mathbf{x}, \sigma) = g_\sigma * u(\mathbf{x}) = \int_{\mathbb{R}^2} g_\sigma(\mathbf{y})u(\mathbf{x} - \mathbf{y})d\mathbf{y}.$$

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

$$217 \quad \frac{\partial}{\partial t}(k_t(\mathbf{x})) = \frac{1}{2}\Delta k_t(\mathbf{x}).$$

218 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
 219 initial condition u (in particular, it is a \mathcal{C}^∞ function on $\mathbb{R}^2 \times (0, \infty)$).

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

$$222 \quad \sigma \frac{\partial g_\sigma}{\partial \sigma}(\mathbf{x}) = \sigma^2 \Delta g_\sigma(\mathbf{x}) = \left(\frac{|\mathbf{x}|^2 - 2\sigma^2}{2\pi\sigma^4}\right) \exp\left(-\frac{|\mathbf{x}|^2}{2\sigma^2}\right).$$

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

$$224 \quad D_u(\mathbf{x}, \sigma) := \sigma^2 \Delta g_\sigma * u(\mathbf{x}) = \sigma^2 \Delta(g_\sigma * u)(\mathbf{x}).$$

225 The following proposition which is recalled without proof shows that these key-
 226 points are covariant to several image transformations.

227 **PROPOSITION 1 ([31]).** *We have the following invariance properties.*

- 228 1. $\forall a \in \mathbb{R}, D_{au} = aD_u$.
- 229 2. *If v is an affine function of \mathbf{x} , then $D_{u+v} = D_u$.*
- 230 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*

$$231 \quad D_{\tau_{\mathbf{h}}u}(\mathbf{x}, \sigma) = D_u(\mathbf{x} - \mathbf{h}, \sigma).$$

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

$$233 \quad D_u(\mathbf{x}, \sigma) = D_v(s\mathbf{x}, s\sigma).$$

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

237 The authors of [46] also discussed the effect of several other image transformations
 238 on the SIFT keypoints but left aside the factor σ^2 in the definition of D_u .

239 **2.2. SIFT Summary.** In the paper by Lowe [33], the scale-normalized LoG is
 240 approximated by a finite difference of Gaussian functions: for a constant scale factor
 241 $k > 1$, he considers instead

$$242 \quad (1) \quad (\mathbf{x}, \sigma) \mapsto (g_{k\sigma} - g_\sigma)(\mathbf{x}) \approx (k\sigma - \sigma) \frac{\partial g_\sigma}{\partial \sigma}(\mathbf{x}) = (k-1)\sigma^2 \Delta g_\sigma(\mathbf{x}).$$

243 Also, the practical implementation of [33] only works with discretized images, so that
 244 the extracted keypoints are actually strict local extrema computed on a discretized
 245 scale-space.

246 Here is a quick summary of the original SIFT method [33]. For technical details
 247 we refer the reader to [54]. Here, and in the remaining of the paper, u_0 refers to the
 248 original image on which we compute keypoints and local descriptors.

249 1. Computing SIFT keypoints:

- 250 (a) Extract local extrema of a discrete version of (1).
- 251 (b) Refine the positions of the local extrema in position and scale using a
 252 quadratic approximation.
- 253 (c) Discard extrema with low contrast (thresholding low values of (1)) and
 254 extrema located on edges (thresholding high values of the ratio between
 255 Hessian eigenvalues).

256 2. Computing SIFT local descriptors associated with the keypoint (\mathbf{x}, σ) :

- 257 (a) Compute one or several principal orientations α . For that, in a square
 258 of size $9\sigma \times 9\sigma$ centered at \mathbf{x} (and parallel to the image axes), compute a
 259 smoothed histogram of orientations of $\nabla g_\sigma * u_0$, and extract its significant
 260 local maxima.
- 261 (b) For each detected orientation α , consider a grid of 4×4 square regions
 262 around (\mathbf{x}, σ) . These square regions, which we call SIFT subcells, are
 263 of size $3\sigma \times 3\sigma$ with one side parallel to α . In each subcell compute the
 264 histogram of $\text{Angle}(\nabla g_\sigma * u_0) - \alpha$ quantized on 8 values ($\ell \frac{\pi}{4}, 1 \leq \ell \leq 8$).
- 265 (c) Normalization: the 16 histograms are concatenated to obtain a feature
 266 vector $f \in \mathbb{R}^{128}$, which is thresholded and normalized

$$267 \quad (2) \quad f_k \leftarrow \min(f_k, 0.2\|f\|_2), \quad f_k = \min\left(255, \left\lfloor 512 \frac{f_k}{\|f\|_2} \right\rfloor\right)$$

268 and finally quantized to 8-bit integers.

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

273 **2.3. Keypoints and Descriptors used in our method.** In the reconstruction
 274 models proposed in this paper, we work with images defined on a rectangle $\Omega \subset \mathbb{Z}^2$ and

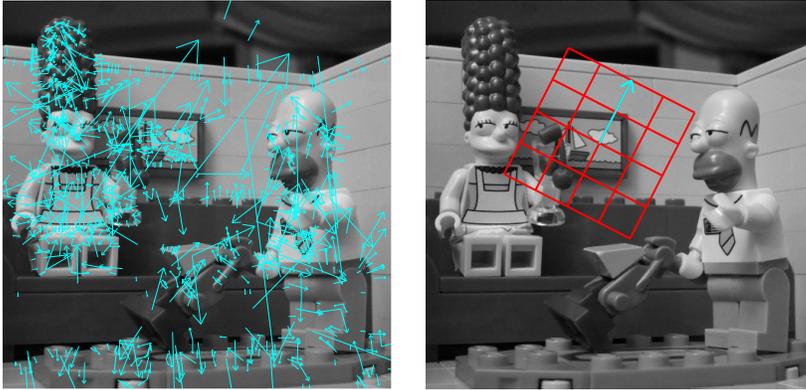


FIG. 1. *Examples of SIFT keypoints and subcells. On the left, one can see an original image (Courtesy of J. Delon) with overimposed SIFT oriented keypoints $(\mathbf{x}, \sigma, \alpha)$ represented as arrows originating from \mathbf{x} , with orientation α and length 6σ . On the right, we display the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

275 we consider the oriented keypoints extracted by the original SIFT method. However,
 276 we will only work with simplified SIFT descriptors in the sense that we extract hard-
 277 binned histograms of gradient orientations at several scales. In other words, we do
 278 not include the vote weights nor the normalization step 2(c).

279 We thus denote by $(s_j)_{j \in \mathcal{J}}$ the collection of SIFT subcells, $s_j \subset \Omega$ (if a $3\sigma \times 3\sigma$
 280 subcell is not entirely contained in Ω , then we replace it with its intersection with Ω).
 281 The SIFT *subcells* must not be confounded with the SIFT *cells*: in a SIFT cell, there
 282 are 16 SIFT subcells so that different subcells s_j can correspond to the same keypoint.
 283 We will denote by $(\mathbf{x}_j, \sigma_j, \alpha_j)$ the oriented keypoint associated with s_j . For $\mathbf{y} \in \Omega$, we
 284 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} .
 285 See Fig. 1 for an illustration.

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

291 **3. Stochastic Models for Gradient Orientations.** In this section, we pro-
 292 pose a model for generating random images constrained to have prescribed local HOGs
 293 in the SIFT subcells. When designing such a model, the main difficulty arises from
 294 the fact that several SIFT subcells can overlap, and thus one has to combine the
 295 information available in all corresponding local HOGs in a way that finally complies
 296 with all the statistical constraints. In order to cope with this issue, we exploit the
 297 framework of exponential distributions to design stochastic orientation models with
 298 prescribed statistical features. The obtained distribution is “as uniform (random) as
 299 possible” in the sense that it is of maximal entropy among all absolutely continuous
 300 distributions which satisfy the desired constraints. We combine this random orienta-
 301 tion field with a deterministic magnitude (which is computed with the scales of locally
 302 available keypoints) in order to obtain a random objective vector field for the gradi-
 303 ent. Finally we solve a Poisson reconstruction problem in order to get back a random
 304 image whose gradient is as close as possible as the randomly sampled objective vector
 305 field.

306 **3.1. Exponential Models with local HOG.** Recall that $\Omega \subset \mathbb{Z}^2$ is a discrete
 307 rectangle. We will denote by $\mathbb{T} = \mathbb{R}/2\pi\mathbb{Z}$ the set of angles, and \mathbb{T}^Ω the set of all
 308 possible orientation fields $\theta = (\theta(\mathbf{x}))_{\mathbf{x} \in \Omega}$ on Ω .

309 **Extracted Statistics.** For simplicity, in contrast with the usual SIFT method,
 310 in this section we only extract gradient orientations at scale 0 and besides we adopt
 311 the same quantization bins for all SIFT subcells

$$312 \quad (3) \quad B_\ell = \left[(\ell - 1) \frac{\pi}{4}, \ell \frac{\pi}{4} \right), \quad (1 \leq \ell \leq 8)$$

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

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

$$316 \quad (4) \quad \forall \theta \in \mathbb{T}^\Omega, \quad f_{j,\ell}(\theta) = \frac{1}{|s_j|} \sum_{\mathbf{x} \in s_j} \mathbf{1}_{B_\ell}(\theta(\mathbf{x})).$$

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

318 **Maximum Entropy Distribution.** We are then interested in probability dis-
 319 tributions P on \mathbb{T}^Ω such that

$$320 \quad (5) \quad \forall j \in \mathcal{J}, \forall \ell \in \{1, \dots, 8\}, \quad \mathbb{E}_P(f_{j,\ell}(\Theta)) = \frac{1}{|s_j|} \sum_{\mathbf{x} \in s_j} \mathbb{P}(\theta(\mathbf{x}) \in B_\ell) = f_{j,\ell}(\theta_0),$$

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

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

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

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

$$336 \quad (6) \quad dP_\lambda = \frac{1}{Z_\lambda} \exp \left(- \sum_{j,\ell} \lambda_{j,\ell} f_{j,\ell}(\theta) \right) d\theta,$$

337 where the partition function Z_λ is given by $Z_\lambda = \int_{\mathbb{T}^\Omega} \exp \left(- \sum_{j,\ell} \lambda_{j,\ell} f_{j,\ell}(\theta) \right) d\theta$, sat-
 338 isfies the constraints (5) and is of maximal entropy among all absolutely continuous
 339 probability distributions w.r.t. the Lebesgue measure $d\theta$ on \mathbb{T}^Ω satisfying the con-
 340 straints (5).

341 *Proof.* This result directly follows from the general theorem given in [47]. The
 342 only difficulty is to handle the hypothesis of linear independence of the $f_{j,\ell}$. In our
 343 framework, the $f_{j,\ell}$ are not independent (in particular because $\sum_{\ell=1}^8 f_{j,\ell} = 1$, and also

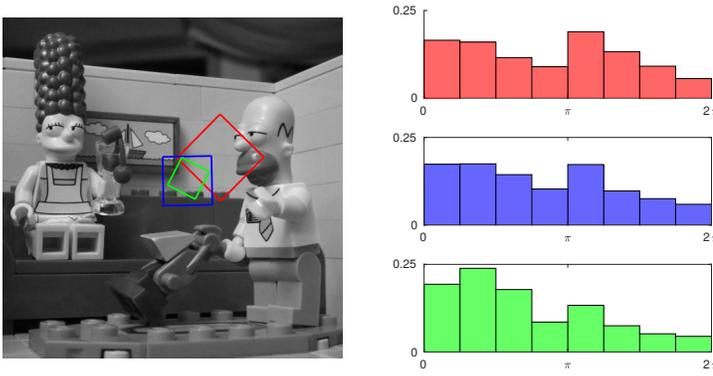


FIG. 2. **Extracting HOG in SIFT subcells.** On the left, we display an original image (Courtesy of J. Delon) with three overlaid SIFT subcells s_j , and on the right, we display the corresponding HOG $(f_{j,\ell}(\theta))_{1 \leq \ell \leq 8}$ extracted in these subcells. The MaxEnt model is a probability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

344 because there may be other dependencies for instance when one subcell is exactly the
 345 union of two smaller subcells). But one can still apply the theorem to an extracted
 346 linearly independent subfamily. This gives the existence of the solution for the initial
 347 family $(f_{j,\ell})$ (but of course not the unicity). \square

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

$$353 \quad (7) \quad D(P||P_\lambda) = \int \log \left(\frac{P(\theta)}{P_\lambda(\theta)} \right) P(\theta) d\theta = -H(P) + \log Z_\lambda + \sum \lambda_{j,\ell} f_{j,\ell}(\theta_0),$$

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

358 Independence Property of the MaxEnt Model.

359 PROPOSITION 3. Under P_λ the values $\Theta(\mathbf{x})$ are independent. Besides, the prob-
 360 ability density function of $\Theta(\mathbf{x})$ is given by

$$361 \quad (8) \quad \frac{1}{Z_{\lambda,\mathbf{x}}} e^{-\varphi_{\lambda,\mathbf{x}}} = \frac{1}{Z_{\lambda,\mathbf{x}}} \sum_{\ell=1}^8 \exp \left(- \sum_{j \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{j,\ell}}{|s_j|} \right) \mathbf{1}_{B_\ell}$$

362

$$363 \quad (9) \quad \text{where } Z_{\lambda,\mathbf{x}} = \sum_{\ell=1}^8 \exp \left(- \sum_{j \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{j,\ell}}{|s_j|} \right) |B_\ell|.$$

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

$$366 \quad (10) \quad -\log \frac{dP_\lambda}{d\theta} - \log Z_\lambda = \sum_{j \in \mathcal{J}, 1 \leq \ell \leq 8} \lambda_{j,\ell} f_{j,\ell}(\theta) = \sum_{\mathbf{x} \in \Omega} \varphi_{\lambda,\mathbf{x}}(\theta(\mathbf{x})),$$

$$367 \quad (11) \quad \text{where} \quad \varphi_{\lambda, \mathbf{x}} = \sum_{\ell=1}^8 \left(\sum_{j \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{j, \ell}}{|s_j|} \right) \mathbf{1}_{B_\ell}.$$

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

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

$$373 \quad (12) \quad \mathbb{E}_{P_\lambda}[\psi(\Theta(\mathbf{x}))] = \frac{\sum_{\ell=1}^8 \exp\left(-\sum_{j \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{j, \ell}}{|s_j|}\right) \int_{B_\ell} \psi(t) dt}{\sum_{\ell=1}^8 \exp\left(-\sum_{j \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{j, \ell}}{|s_j|}\right) |B_\ell|}$$

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

$$376 \quad (13) \quad \mathbb{E}_{P_\lambda}[f_{j, \ell}(\Theta)] = \frac{1}{|s_j|} \sum_{\mathbf{x} \in s_j} \mathbb{P}(\Theta(\mathbf{x}) \in B_\ell) = \frac{1}{|s_j|} \sum_{\mathbf{x} \in s_j} \frac{\exp\left(-\sum_{k \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{k, \ell}}{|s_k|}\right) |B_\ell|}{\sum_{1 \leq \ell' \leq 8} \exp\left(-\sum_{k \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{k, \ell'}}{|s_k|}\right) |B_{\ell'}|}.$$

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

$$379 \quad (14) \quad \forall j, \ell, \quad \sum_{\mathbf{x} \in s_j} \frac{1}{Z_{\lambda, \mathbf{x}}} \exp\left(-\sum_{k \in \mathcal{J}(\mathbf{x})} \frac{\lambda_{k, \ell}}{|s_k|}\right) |B_\ell| = |\{\mathbf{x} \in s_j ; \theta_0(\mathbf{x}) \in B_\ell\}|.$$

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

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

$$383 \quad (15) \quad Z_{\lambda, \mathbf{x}} = \sum_{\ell=1}^8 \exp\left(-\frac{\lambda_{j, \ell}}{|s_j|}\right) |B_\ell|.$$

384 Then (14) gives

$$385 \quad (16) \quad \forall \ell, \quad \frac{1}{Z_{\lambda, \mathbf{x}}} \exp\left(-\frac{\lambda_{j, \ell}}{|s_j|}\right) = \frac{|\{\mathbf{x} \in s_j ; \theta_0(\mathbf{x}) \in B_\ell\}|}{|s_j| |B_\ell|} = f_{j, \ell}(\theta_0),$$

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

$$387 \quad (17) \quad \frac{1}{Z_{\lambda, \mathbf{x}}} e^{-\varphi_{\lambda, \mathbf{x}}} = \sum_{\ell=1}^8 \frac{|\{\mathbf{x} \in s_j ; \theta_0(\mathbf{x}) \in B_\ell\}|}{|s_j| |B_\ell|} \mathbf{1}_{B_\ell} = \sum_{\ell=1}^8 f_{j, \ell}(\theta_0) \frac{1}{|B_\ell|} \mathbf{1}_{B_\ell}.$$

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

$$393 \quad (18) \quad \forall \ell, \quad \lambda_{j, \ell} = -|s_j|(\log f_{j, \ell}(\theta_0) + \log c_j).$$

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

$$399 \quad (19) \quad \Phi(\lambda) = \log Z_\lambda + \sum_{j,\ell} \lambda_{j,\ell} f_{j,\ell}(\theta_0).$$

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

$$401 \quad (20) \quad \frac{\partial \log Z_\lambda}{\partial \lambda_{j,\ell}} = \frac{1}{Z_\lambda} \frac{\partial Z_\lambda}{\partial \lambda_{j,\ell}} = -\mathbb{E}_{P_\lambda} [f_{j,\ell}(\Theta)],$$

402 which gives

$$403 \quad (21) \quad \frac{\partial \Phi}{\partial \lambda_{j,\ell}} = f_{j,\ell}(\theta_0) - \mathbb{E}_{P_\lambda} [f_{j,\ell}(\Theta)].$$

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

405 Similarly, we can also obtain the second order derivatives

$$406 \quad (22) \quad \frac{\partial^2 \Phi}{\partial \lambda_{j,\ell} \partial \lambda_{j',\ell'}} = \mathbb{E}_{P_\lambda} \left[(f_{j,\ell}(\Theta) - \mathbb{E}_{P_\lambda} [f_{j,\ell}(\Theta)]) (f_{j',\ell'}(\Theta) - \mathbb{E}_{P_\lambda} [f_{j',\ell'}(\Theta)]) \right].$$

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

411 Therefore, we can compute the solution P_λ by a gradient descent algorithm in
 412 order to minimize Φ . The complete algorithm is summarized in Section 3.3. Since Φ is
 413 not strictly convex, we will not have a guarantee of convergence on the iterates, but on
 414 the function values. Since $|f_{j,\ell}(\theta)| \leq 1$, it is straightforward to see that all coefficients
 415 of the Hessian $\nabla^2 \Phi(\lambda)$ have modulus ≤ 1 . Therefore, the ℓ^2 operator norm of $\nabla^2 \Phi$
 416 is bounded by $8|\mathcal{J}|$, which implies that $\nabla \Phi$ is L -Lipschitz with $L = 8|\mathcal{J}|$. Writing
 417 λ^k the iterates of the gradient descent with constant step size $h < \frac{2}{L}$, [49, Th 2.1.14]
 418 gives

$$419 \quad (23) \quad \Phi(\lambda^k) - \min \Phi = \mathcal{O}\left(\frac{1}{k}\right).$$

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

424 **3.2. Monoscale Poisson Reconstruction.** Now that we have built a random
 425 orientation field Θ with maximum entropy distribution P_λ , we will use it to propose
 426 a target vector field V for the image gradient. More precisely, we set the gradient
 427 magnitude at \mathbf{x} in a deterministic manner, as the inverse scale of the smallest subcell
 428 that covers \mathbf{x} . For pixels \mathbf{x} which lie outside the SIFT subcells, we set $V(\mathbf{x}) = 0$.
 429 This choice allows to give more weight to the locations for which we have information
 430 at finer scale. It is also motivated by the following homogeneity argument. Assume

431 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
 432 Proposition 1, v has a keypoint $(a\mathbf{x}, a\sigma)$. Let us compare the mean gradient magnitude
 433 at scale σ in the corresponding subcell s to the analogous quantity for v . A simple
 434 computation shows that

$$435 \quad \frac{1}{|as|} \int_{\lambda s} |\nabla g_{a\sigma} * v(\mathbf{y})| d\mathbf{y} = \frac{1}{a} \frac{1}{|s|} \int_s |\nabla g_{\sigma} * u(\mathbf{y})| d\mathbf{y},$$

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

$$444 \quad (24) \quad \forall \mathbf{x} \in \Omega, \quad V(\mathbf{x}) = \left(\max_{j \in \mathcal{J}(\mathbf{x})} \frac{1}{\sigma_j} \right) e^{i\Theta(\mathbf{x})} \mathbf{1}_{\mathcal{J}(\mathbf{x}) \neq \emptyset}.$$

445 The aim of the Poisson reconstruction is to compute an image whose gradient is
 446 as close as possible to the vector field $V = (V_1, V_2)$. In the case of image editing,
 447 this technique has been proposed by Pérez et al. [52] in order to copy pieces of an
 448 image into another one in a seamless way. More precisely, the goal is to minimize the
 449 functional

$$450 \quad (25) \quad F(u) = \sum_{\mathbf{x} \in \Omega} \|\nabla u(\mathbf{x}) - V(\mathbf{x})\|_2^2.$$

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

$$453 \quad (26) \quad U = \text{Argmin}\{F(u); u : \Omega \rightarrow \mathbb{R} \text{ and such that } \sum_{\mathbf{x} \in \Omega} u(\mathbf{x}) = 0\}.$$

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

$$457 \quad (27) \quad \nabla u(\mathbf{x}) = \begin{pmatrix} \partial_1 * u(\mathbf{x}) \\ \partial_2 * u(\mathbf{x}) \end{pmatrix} \quad \text{where} \quad \begin{cases} \partial_1 &= \delta_{(0,0)} - \delta_{(1,0)} \\ \partial_2 &= \delta_{(0,0)} - \delta_{(0,1)} \end{cases},$$

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

$$459 \quad (28) \quad F(u) = \frac{1}{|\Omega|} \sum_{\boldsymbol{\xi}} |\widehat{\partial}_1(\boldsymbol{\xi}) \widehat{u}(\boldsymbol{\xi}) - \widehat{V}_1(\boldsymbol{\xi})|_2^2 + |\widehat{\partial}_2(\boldsymbol{\xi}) \widehat{u}(\boldsymbol{\xi}) - \widehat{V}_2(\boldsymbol{\xi})|_2^2.$$

460
 461 Thus, for each $\boldsymbol{\xi}$ we have a barycenter problem which is simply solved by

$$462 \quad (29) \quad \forall \boldsymbol{\xi} \neq 0, \quad \widehat{U}(\boldsymbol{\xi}) = \frac{\overline{\widehat{\partial}_1(\boldsymbol{\xi})} \widehat{V}_1(\boldsymbol{\xi}) + \overline{\widehat{\partial}_2(\boldsymbol{\xi})} \widehat{V}_2(\boldsymbol{\xi})}{|\widehat{\partial}_1(\boldsymbol{\xi})|^2 + |\widehat{\partial}_2(\boldsymbol{\xi})|^2} \quad \text{and} \quad \widehat{U}(0) = 0.$$

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

465 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

466 (30)
$$\widehat{U}(\boldsymbol{\xi}) = \widehat{\nu}(\boldsymbol{\xi})\widehat{V}(\boldsymbol{\xi}) \quad \text{where} \quad \widehat{\nu}(\boldsymbol{\xi}) = \begin{cases} \frac{\widehat{\nabla}(\boldsymbol{\xi})^*}{|\widehat{\nabla}(\boldsymbol{\xi})|^2} & \text{if } \boldsymbol{\xi} \neq 0 \\ 0 & \text{if } \boldsymbol{\xi} = 0 \end{cases}.$$

467 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

468 (31)
$$U = \nu * V = \nu_1 * V_1 + \nu_2 * V_2.$$

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

472 **3.3. Algorithm.** Here we summarize the algorithm for estimating and sampling
473 the MaxEnt model proposed in this section. In Fig. 3 we display an example of
474 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)
$$V(\boldsymbol{x}) = \left(\max_{j \in \mathcal{J}(\boldsymbol{x})} \frac{1}{\sigma_j} \right) e^{i\theta(\boldsymbol{x})} \mathbf{1}_{\mathcal{J}(\boldsymbol{x}) \neq \emptyset}$$

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

475 For images having many SIFT keypoints in overlapping positions, this algorithm
476 may be slow to converge as can be observed on the case of Fig. 3. This case is
477 relatively simple because it has only 187 keypoints but this corresponds already to
478 $8 \times 16 \times 187 \approx 24000$ $\lambda_{j,\ell}$ parameters to estimate. This is why we use a stopping
479 criterion based on a maximal number of iterations.

480 **3.4. Discussion on MaxEnt Model.** One drawback of MaxEnt is that the
481 guarantee on the local distributions of orientations is lost after the Poisson recon-
482 struction step. One way to cope with that would be to consider a model that operates
483 directly on the image values, and not on the orientation field. Theorem 2 could be
484 extended to statistics like

485 (33)
$$\tilde{f}_{j,\ell}(u) = \frac{1}{|s_j|} \sum_{\boldsymbol{x} \in s_j} \mathbf{1}_{B_\ell}(\text{Angle}(\nabla u(\boldsymbol{x}))).$$

486 It is even possible to consider multiscale statistics using $\nabla g_{\sigma_j} * u$ instead of ∇u (as
487 it will be the case in Section 4). But the analog of Proposition 3 would not hold for
488 these models, so that sampling should rely on a Gibbs strategy. Its cost would be

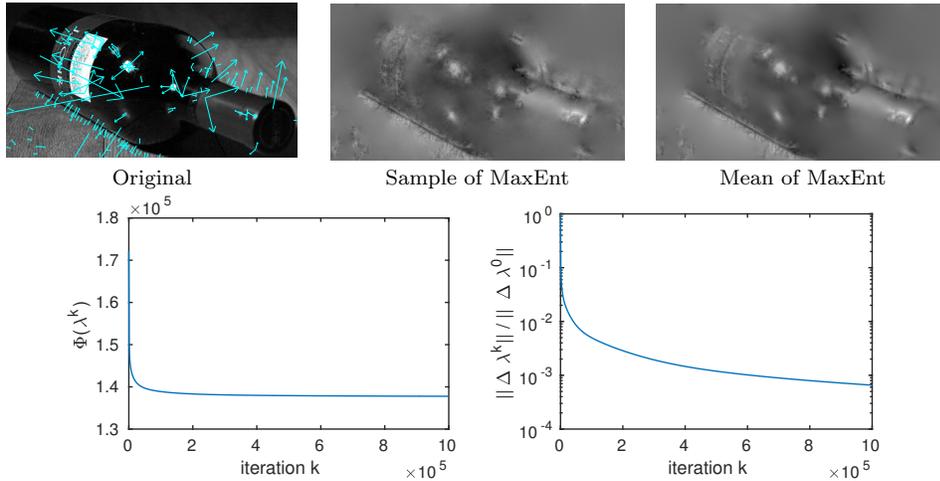


FIG. 3. **Reconstruction with the MaxEnt model.** In the first row from left to right, we display an original image with overimposed 187 oriented keypoints, a sample of the associated MaxEnt model, and the expectation of the MaxEnt model. In the second row we display the evolution of Φ along the iterates,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

489 clearly prohibitive in the multiscale case due to the large Markov neighborhood size.
 490 Even in the monoscale case the convergence of this Gibbs sampler may be very long
 491 depending on the parameters λ ; and since we would need one sample per iteration of
 492 gradient descent to estimate λ , we chose to leave it aside and concentrate on models
 493 with reasonably fast sampling.

494 Also, one can consider another orientation model in which the local HOGs are
 495 computed with a quantization that depends on the keypoint orientation. The inde-
 496 pendence property still holds for this model, and the marginal orientations still have
 497 a piecewise constant density, but the number of parameters would be much larger
 498 (there would be as many ℓ 's as bins of a subdivision that is adapted to all keypoint
 499 orientations). Therefore this model is practically untractable, and also only of minor
 500 interest. Indeed, in view of the results of Fig. 3, it is likely that the used quantization
 501 has only a minor impact on the visual results (provided that we still have a minimal
 502 number of bins).

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

511 **4.1. Construction of MS-Poisson Model.**

512 **Extracted Statistics.** The MS-Poisson model is based on local statistics on
 513 multiscale gradient orientations. More precisely, in s_j we extract the quantized HOG

514 at scale σ_j

$$515 \quad (34) \quad H_{j,\ell} = \frac{1}{|s_j|} \left| \left\{ \mathbf{x} \in s_j ; \text{Angle}(\nabla g_{\sigma_j} * u_0)(\mathbf{x}) - \alpha_j \in [(\ell-1)\frac{\pi}{4}, \ell\frac{\pi}{4}] \right\} \right|.$$

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

$$518 \quad (35) \quad h_j = \frac{4}{\pi} \sum_{\ell=1}^8 H_{j,\ell} \mathbf{1}_{[\alpha_j + (\ell-1)\frac{\pi}{4}, \alpha_j + \ell\frac{\pi}{4}]}.$$

519 Notice that, in contrast to the statistics (4) used in the MaxEnt model, the quanti-
520 zation here depends on the local orientation α_j .

521 **Target Vector Fields at Multiple Scales.** Using the local orientation distri-
522 butions h_j , we define vector fields $V_j : \Omega \rightarrow \mathbb{R}^2$ that will serve as objective gradients
523 at scale σ_j in the SIFT subcell s_j . We propose to set

$$524 \quad (36) \quad \forall \mathbf{x} \in \Omega, \quad V_j(\mathbf{x}) = \frac{1}{\sigma_j} e^{i\gamma_j(\mathbf{x})} \mathbf{1}_{s_j}(\mathbf{x}),$$

525 where the orientations $\gamma_j(\mathbf{x})$ are independently sampled according to the distribu-
526 tion h_j . Again, as justified in Section 3.2, we set the gradient magnitude in a deter-
527 ministic way using the inverse of the scale σ_j . Once these vector fields V_j have been
528 sampled, we obtain an image U by solving a multiscale Poisson problem as explained
529 in the next paragraph.

530 **4.2. Multiscale Poisson Reconstruction.** In order to simultaneously con-
531 strain the gradient at several scales $(\sigma_j)_{j \in \mathcal{J}}$, we propose to consider the following
532 multiscale Poisson energy

$$533 \quad (37) \quad G(u) = \sum_{j \in \mathcal{J}} w(\sigma_j) \sum_{\mathbf{x} \in \Omega} \|\nabla(g_{\sigma_j} * u)(\mathbf{x}) - V_j(\mathbf{x})\|_2^2,$$

534 where g_σ is the Gaussian kernel of standard deviation σ , $V_j = (V_{j,1}, V_{j,2})^T$ is the
535 objective gradient at scale σ_j , and $\{w(\sigma_j), j \in \mathcal{J}\}$ is a set of weights. In our applica-
536 tion, since there are more keypoints in the fine scales (i.e. with small σ_j), and since
537 the keypoints at fine scales are generally more informative, a reasonable choice is to
538 take all weights $w(\sigma_j) = 1$. But we keep these weights in the formula for the sake of
539 generality. We thus set

$$540 \quad (38) \quad U = \text{Argmin}\{G(u); u : \Omega \rightarrow \mathbb{R} \text{ and such that } \sum_{\mathbf{x} \in \Omega} u(\mathbf{x}) = 0\}.$$

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

$$543 \quad (39) \quad G(u) = \frac{1}{|\Omega|} \sum_{j \in \mathcal{J}} \sum_{\boldsymbol{\xi}} w(\sigma_j) \left(|\widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \widehat{\partial}_1(\boldsymbol{\xi}) \widehat{u}(\boldsymbol{\xi}) - \widehat{V}_{j,1}(\boldsymbol{\xi})|_2^2 + |\widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \widehat{\partial}_2(\boldsymbol{\xi}) \widehat{u}(\boldsymbol{\xi}) - \widehat{V}_{j,2}(\boldsymbol{\xi})|_2^2 \right).$$

544 As for the monoscale Poisson problem, the solution U is still a barycenter given by
545 $\widehat{U}(0) = 0$ and

$$546 \quad (40) \quad \forall \boldsymbol{\xi} \neq 0, \quad \widehat{U}(\boldsymbol{\xi}) = \frac{\sum_{j \in \mathcal{J}} w(\sigma_j) \widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \left(\overline{\widehat{\partial}_1(\boldsymbol{\xi}) \widehat{V}_{j,1}(\boldsymbol{\xi})} + \overline{\widehat{\partial}_2(\boldsymbol{\xi}) \widehat{V}_{j,2}(\boldsymbol{\xi})} \right)}{\sum_{j \in \mathcal{J}} w(\sigma_j) |\widehat{g}_{\sigma_j}(\boldsymbol{\xi})|^2 \left(|\widehat{\partial}_1(\boldsymbol{\xi})|^2 + |\widehat{\partial}_2(\boldsymbol{\xi})|^2 \right)}.$$

547 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.

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

$$552 \quad (41) \quad U = \operatorname{Argmin}\{G(u) + \mu \|\nabla u\|_2^2; u : \Omega \rightarrow \mathbb{R} \text{ and such that } \sum_{\mathbf{x} \in \Omega} u(\mathbf{x}) = 0\},$$

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

$$554 \quad (42) \quad \forall \boldsymbol{\xi} \neq 0, \quad \widehat{U}(\boldsymbol{\xi}) = \frac{\sum_{j \in \mathcal{J}} w(\sigma_j) \widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \left(\overline{\widehat{\partial}_1(\boldsymbol{\xi})} \widehat{V}_{j,1}(\boldsymbol{\xi}) + \overline{\widehat{\partial}_2(\boldsymbol{\xi})} \widehat{V}_{j,2}(\boldsymbol{\xi}) \right)}{\left(\mu + \sum_{j \in \mathcal{J}} w(\sigma_j) |\widehat{g}_{\sigma_j}(\boldsymbol{\xi})|^2 \right) \left(|\widehat{\partial}_1(\boldsymbol{\xi})|^2 + |\widehat{\partial}_2(\boldsymbol{\xi})|^2 \right)}.$$

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

559 We end this paragraph 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 \widehat{U} by solving the MS-Poisson problem (41) with targets V_j , with $w(\sigma_j) = 1$ and $\mu = 50$.
-

560 **Remark:** In (42), one can observe that the solution to MS-Poisson actually solves a
 561 monoscale Poisson problem with objective vector field V whose Fourier transform is
 562 given by

$$563 \quad (43) \quad \widehat{V}(\boldsymbol{\xi}) = \frac{\sum_{j \in \mathcal{J}} w(\sigma_j) \widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \widehat{V}_j(\boldsymbol{\xi})}{\mu + \sum_{j \in \mathcal{J}} w(\sigma_j) |\widehat{g}_{\sigma_j}(\boldsymbol{\xi})|^2}.$$

564 **4.3. First and Second Order Moments.** In order to compute the statistics
 565 of the MS-Poisson model, we remark that the multiscale Poisson reconstruction is
 566 actually a linear process. Indeed, for each j , let $\nu_j : \Omega \rightarrow \mathbb{R}^{1 \times 2}$ be the vector-valued
 567 kernel defined by its discrete Fourier transform

$$568 \quad (44) \quad \forall \boldsymbol{\xi} \neq 0, \quad \widehat{\nu}_j(\boldsymbol{\xi}) = \frac{w(\sigma_j) \widehat{g}_{\sigma_j}(\boldsymbol{\xi}) \widehat{\nabla}(\boldsymbol{\xi})^*}{\left(\mu + \sum_{j' \in \mathcal{J}} w(\sigma_{j'}) |\widehat{g}_{\sigma_{j'}}(\boldsymbol{\xi})|^2 \right) |\widehat{\nabla}(\boldsymbol{\xi})|^2} \quad \text{and } \widehat{\nu}_j(0) = 0.$$

569 Then, as in Section 3.2 we get the convolutive expression

$$570 \quad (45) \quad U = \sum_{j \in \mathcal{J}} \nu_j * V_j = \sum_{j \in \mathcal{J}} \left(\nu_{j,1} * V_{j,1} + \nu_{j,2} * V_{j,2} \right).$$

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

$$572 \quad (46) \quad \mathbb{E}(U) = \sum_{j \in \mathcal{J}} \nu_j * \mathbb{E}(V_j),$$

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

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

$$576 \quad (47) \quad \text{Var}(U(\mathbf{x})) = \sum_{j \in \mathcal{J}} \text{Var}(\nu_j * V_j(\mathbf{x})).$$

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

$$578 \quad (48) \quad \text{Var}(\nu_j * V_j(\mathbf{x})) = \text{Var}\left(\sum_{\mathbf{y} \in \Omega} \nu_j(\mathbf{x} - \mathbf{y}) V_j(\mathbf{y})\right) = \sum_{\mathbf{y} \in \Omega} \text{Var}(\nu_j(\mathbf{x} - \mathbf{y}) V_j(\mathbf{y}))$$

$$579 \quad (49) \quad = \sum_{\mathbf{y} \in \Omega} \nu_j(\mathbf{x} - \mathbf{y}) \text{Cov}(V_j(\mathbf{y})) \nu_j^T(\mathbf{x} - \mathbf{y})$$

$$580 \quad (50) \quad = \sum_{\mathbf{y} \in \Omega} \nu_{j,1}^2(\mathbf{x} - \mathbf{y}) \text{Var}(V_{j,1}(\mathbf{y})) + \nu_{j,2}^2(\mathbf{x} - \mathbf{y}) \text{Var}(V_{j,2}(\mathbf{y}))$$

$$581 \quad (51) \quad + 2\nu_{j,1}(\mathbf{x} - \mathbf{y}) \nu_{j,2}(\mathbf{x} - \mathbf{y}) \text{Cov}(V_{j,1}(\mathbf{y}), V_{j,2}(\mathbf{y})).$$

583 Therefore the variance of this model can be obtained by summing convolutions of the
584 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
585 p.d.f. h_j given by (34), we can explicitly compute its covariance.

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

$$588 \quad (52) \quad \text{Cov}(U(\mathbf{x}), U(\mathbf{y})) = \sum_{j \in \mathcal{J}} \sum_{\mathbf{z} \in \Omega} \nu_j(\mathbf{x} - \mathbf{z}) \text{Cov}(V_j(\mathbf{z})) \nu_j^T(\mathbf{y} - \mathbf{z}).$$

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

599 **5.1. Results with MaxEnt and MS-Poisson model.**

600 **5.1.1. Comparison between MaxEnt and MS-Poisson.** Let us first com-
 601 pare the reconstruction results obtained with MaxEnt and with MS-Poisson. On
 602 Fig. 4, using an original image with 386 keypoints, we display a sample of MaxEnt
 603 and a sample of MS-Poisson, together with the expected images of these models. One
 604 first remark is that both models are able to retrieve several geometric structures of the
 605 original image, so that much semantic content of the image can still be understood.
 606 For both models, one can observe that the samples are very close to the expected
 607 image, which will be later confirmed by the variance analysis on Fig. 8.

608 One crucial difference between MaxEnt and MS-Poisson is that they do not rely
 609 on the same gradient information. Indeed, MS-Poisson exploits gradients extracted
 610 at multiple scales while MaxEnt only operates with gradients at scale $\sigma = 0$ (i.e.
 611 the same scale as the image). This is why the results obtained with MS-Poisson
 612 will generally look blurrier than the ones obtained with MaxEnt. Besides, because
 613 of the multiscale nature of the input of MS-Poisson, the corresponding optimization
 614 problem had to be regularized; and the adopted H^1 -regularization term is also a source
 615 of blur in the result. This is confirmed by Fig. 5 where we display several MS-Poisson
 616 reconstructions with varying regularization parameter μ . In Fig. 5 and in many other
 617 experiments, we observed that the parameter $\mu = 50$ realizes a good compromise
 618 between preserving geometric structures and removing spurious oscillations.

619 In the last row of Fig. 4, we also compare with the reconstructions obtained with
 620 the true gradient orientations (resp. multiscale gradient orientations) computed in
 621 the SIFT subcells and the gradient magnitude computed as in MaxEnt (resp. MS-
 622 Poisson). So the difference with MaxEnt (or MS-Poisson) is that local (multiscale)
 623 gradient orientations are not pooled in histograms but directly extracted pixelwise;
 624 in other words, there is no local resampling of the orientations. Thus, in some sense,
 625 these images are the best ones we could hope using Poisson reconstruction. Comparing
 626 these images with samples of MS-Poisson and MaxEnt precisely shows the effect of
 627 local resampling of the (multiscale) orientations; observe in particular the man's face
 628 and also the folds of his t-shirt. These images thus correspond to much more precise
 629 reconstructions, but it is interesting to notice that in certain regions where attention
 630 will be focused (near the face e.g.), there are enough keypoints at fine scales in order
 631 to get back satisfying pieces of images even after local resampling. Also, one must
 632 keep in mind that the loss of the gradient magnitude information is in practice difficult
 633 to cope with and may force us to erroneously amplify the noise in the reconstruction.
 634 As one can see in the bottom left of Fig. 4, it is obvious if one tries to set the gradient
 635 magnitude to 1 in the global Poisson reconstruction.

636 **5.1.2. Quantitative evaluation.** As mentioned in [15], there is no reliable cri-
 637 terion to quantitatively evaluate the quality of the result for such reconstruction prob-
 638 lems. In our context where only gradient orientations are extracted, it is reasonable
 639 to evaluate the reconstruction quality based on the normalized correlation to the in-
 640 put image (which is invariant under affine contrast change). If $u, v : \Omega \rightarrow \mathbb{R}$ are two
 641 images, the normalized correlation is defined as

$$642 \quad (53) \quad r(u, v) = \frac{1}{|\Omega|} \sum_{\mathbf{x} \in \Omega} \left(\frac{u(\mathbf{x}) - \bar{u}}{\sigma_u} \right) \left(\frac{v(\mathbf{x}) - \bar{v}}{\sigma_v} \right) \in [-1, 1],$$

643 where $\bar{u} = \frac{1}{|\Omega|} \sum_{\mathbf{x} \in \Omega} u(\mathbf{x})$ and $\sigma_u^2 = \frac{1}{|\Omega|} \sum_{\mathbf{x} \in \Omega} (u(\mathbf{x}) - \bar{u})^2$. On Fig. 4, for each re-
 644 sult we have indicated the normalized correlation value r . Surprisingly, the higher
 645 correlation values are attained with results linked to the MS-Poisson model (even if
 646 it only have access to HOG computed on a blurred gradient). Besides, the value



FIG. 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 on the electronic version)

647 attained by the samples (or mean) of MS-Poisson is close to the one obtained with
 648 the true multiscale HOGs. In contrast, the correlations obtained with the MaxEnt
 649 model are lower. This is better explained by the results of Fig. 6 in which we display
 650 values of local normalized correlations obtained with both models: for each pixel \mathbf{x}
 651 we extract patches $p_{\mathbf{x}}(u), p_{\mathbf{x}}(v)$ of compared images u, v and we compute the nor-
 652 malized correlation $r(p_{\mathbf{x}}(u), p_{\mathbf{x}}(v))$. On the one hand, MaxEnt result is everywhere
 653 much noisier (because gradient orientations are sampled independently). On the other
 654 hand, the regularization involved in MS-Poisson helps to propagate good correla-
 655 tions values in regions located near SIFT subcells. Also, this criterion based on normalized



FIG. 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 on the electronic version)*

656 correlation confirms the choice for the regularization parameter $\mu = 50$, see Fig. 5.

657 Another interesting way of performing quantitative evaluation in our context is
 658 to compare the HOG computed in the SIFT subcells to the ones of the original image.
 659 For each subcell, we can compute histograms H_u, H_v of gradient orientations (with 8
 660 bins) for the images u, v and then compute the total variation distance between these
 661 histograms, defined as $\frac{1}{2} \sum_{\ell=1}^8 |H_u(\ell) - H_v(\ell)| \in [0, 1]$. Again, we use gradients at
 662 scale 0 when considering the MaxEnt model, and scaled gradients when considering
 663 the MS-Poisson model. We can then average the HOG distances obtained for all SIFT
 664 subcells, weighted by the number of pixels in each subcell. With this methodology, for
 665 the image of Fig. 4, we obtain a mean distance around 0.27 for MS-Poisson and 0.16
 666 for MaxEnt. This value is lower for MaxEnt because the model is inherently made to
 667 satisfy the HOG constraint. One can better understand these results by examining the
 668 orientation fields of both models, as proposed in Fig. 7, in particular the effect of the
 669 final Poisson reconstruction (keeping in mind that MS-Poisson can also be written
 670 with a single objective vector field given by Eq. (43)). In this figure, one clearly
 671 observes that the objective vector field for MS-Poisson is already very smooth (and
 672 certainly too smooth to account for fine local variations in orientation). In contrast
 673 the objective vector field for MaxEnt better accounts for the fine variations, but is
 674 much noisier, even after the Poisson reconstruction step.

675 **5.1.3. Second order statistics.** As we have seen in Section 4.3, it is possible
 676 to compute the second order statistics of the reconstructed image in each model. In
 677 Fig. 8 we display the standard deviations of all pixels values in each model. One first
 678 remark is that MaxEnt has in general much larger variance than MS-Poisson which can
 679 be explained by the fact that the output of MS-Poisson is in some sense a weighted

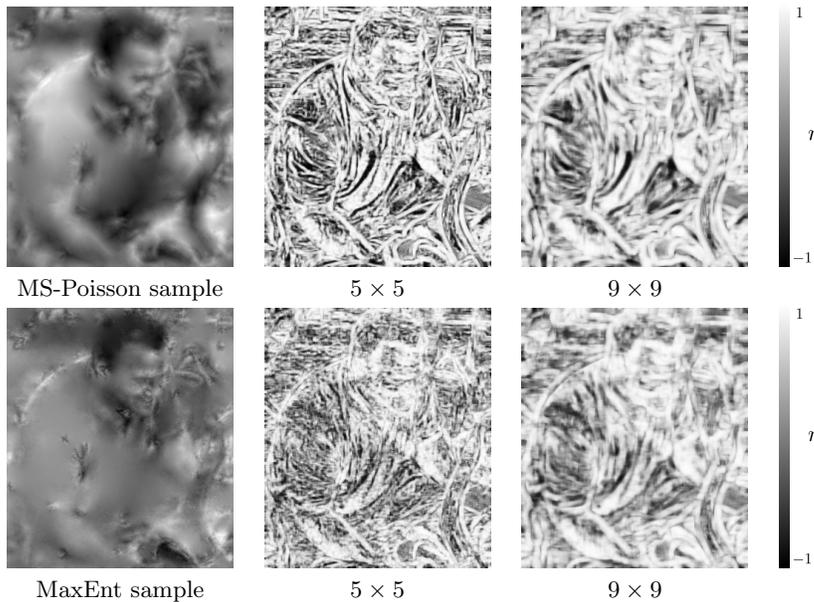


FIG. 6. *Comparison between MS-Poisson and MaxEnt with local normalized correlations.* On the left column we display samples of the models MS-Poisson and MaxEnt. On the other columns, we display the local normalized correlation of the sample (first column) with respect to the original image. The local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

680 average of many local reconstructions. Also it is interesting to see that the image
 681 regions with larger variance are located in the SIFT subcells which contain sharp
 682 geometric details. That being said, the variance of both these models is relatively
 683 small compared to the global range of the mean image, which indicates that both
 684 these models have quite small variations around the mean.

685 **5.1.4. Discard boundary keypoints.** Let us emphasize that in our experi-
 686 ments, we used all the keypoints computed by the SIFT methods and we did not
 687 discard keypoints located near the image boundaries. The positions of the corre-
 688 sponding local extrema in the normalized scale-space are indeed highly dependent on
 689 the boundary conditions used to compute the scale-space. This explains why SIFT
 690 keypoints near the image boundaries are often discarded for particular applications,
 691 e.g. image matching. In our reconstruction problem, there is no reason to discard such
 692 keypoints, and we use the information available in SIFT subcells as soon as they in-
 693 tersect the image domain (if the SIFT subcell is not entirely contained in the domain,
 694 we consider only the pixels in the intersection of the subcell and the domain). But
 695 still, it is clear that for some images, the reconstruction will be quite different when
 696 discarding those keypoints. For example in the case of Fig. 9, if boundary keypoints
 697 are discarded, then several parts of the man’s body are not as properly retrieved in
 698 the reconstruction, thus affecting the semantic understanding of the image.

699 **5.1.5. Matching keypoints between the original and reconstructed im-**
 700 **ages.** Finally, it is interesting to compare the keypoints computed on the original
 701 image and the ones computed on several samples of the models. As one can see on
 702 Fig. 10, we get back similar keypoints in many regions of the image, but still with some

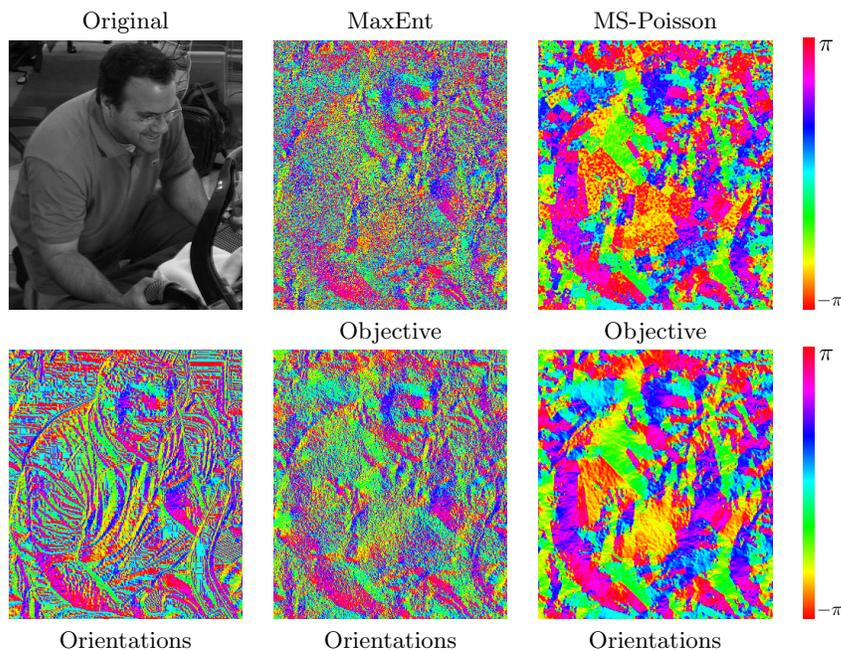


FIG. 7. **Orientation fields of the original and reconstructions.** On the left column we display the original image (top) and the corresponding gradient orientations (bottom). On column 2 (MaxEnt) and column 3 (MS-Poisson), we display the orientation of the objective vector field (before Poisson reconstruction, 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 smoothes slightly the orientation field. (Images are better seen on the electronic version)

703 variations in positions, scales and orientations. In particular, we observe variations
 704 when taking different samples of the model (sometimes, some keypoints associated
 705 with low contrast regions may even disappear). Notice also that we get back less
 706 keypoints in the MS-Poisson model: indeed, since it is more regular we loose some
 707 extrema in the scale-space. Besides, the regularization tends to change the scale of
 708 the structures, thus the scales of the keypoints is often larger than in the original
 709 image.

710 In order to give a more quantitative evaluation of the variations of the keypoints
 711 over different samples of the model, it is possible to use the matching algorithm avail-
 712 able with the online implementation [54] (we used the proposed default parameters).
 713 This algorithm follows the matching method proposed in [33] which essentially pairs
 714 SIFT keypoints by thresholding the ratio between the distances to the first and second
 715 nearest neighbors (computed with the ℓ^2 -distance between SIFT descriptors). First
 716 we can comment on what happens when matching two different samples of the same
 717 model. For the MS-Poisson model, when matching the two samples shown in Fig. 10,
 718 among the 206 keypoints found on the first image (resp. 211 on the second image), 150
 719 keypoints are matched. The mean spatial distance (resp. mean scale variation, mean
 720 angle variation) between matched keypoints is about 0.54 (resp. 0.15, 0.050). Similar
 721 numbers can be given for the MaxEnt model, but in this case much less keypoints are
 722 correctly matched: over the 452 keypoints found on the first image (resp. 458 on the

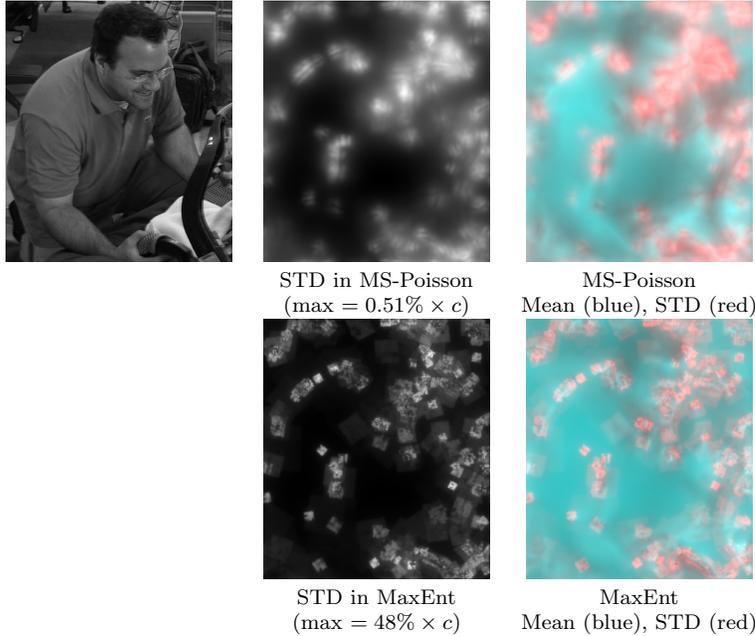


FIG. 8. *Standard deviations of MS-Poisson and MaxEnt models.* On the top left we display the original image. On the rest of the figure we display the images formed with the standard deviations (STD) of the models MS-Poisson (first row) and MaxEnt (second row). On the second column we display the raw 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)

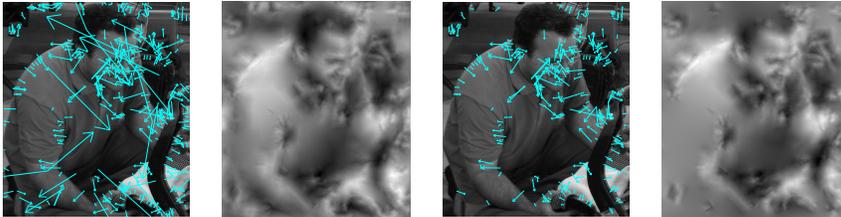


FIG. 9. *Discard keypoints near image boundary.* In this figure, we examine the effect of discarding keypoints whose associated SIFT cell is not entirely contained in the image domain. The displayed reconstructions are samples of the MS-Poisson model.

723 second image), only 184 are matched. This reflects again the larger variance of the
 724 MaxEnt model.

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

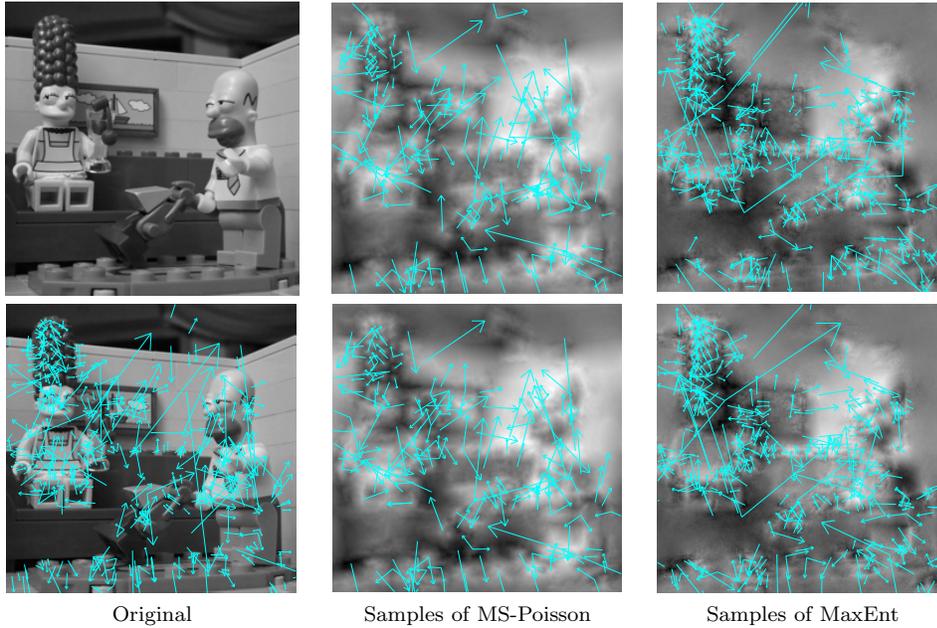


FIG. 10. **Keypoints after reconstruction.** In the first column we display an original image and the same image with its SIFT keypoints. In the second column we display two samples of the MS-Poisson model. In the third column we display two samples of the MaxEnt model. We display the keypoints associated 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)

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

735 **5.2. Reconstruction from true SIFT descriptors.** The two models MS-
 736 Poisson and MaxEnt are designed to propose stochastic reconstructions of an image
 737 based on simplified SIFT descriptors, that is, multiscale HOGs extracted around the
 738 SIFT keypoints. But it is also possible to test these reconstruction models with the
 739 true SIFT descriptors. For that, for each keypoint, we still consider the location, scale
 740 and principal orientation, but, following the discussion of Section 2.2, starting from
 741 the normalized feature vector $(f_k) \in \mathbb{R}^{128}$, we improperly build target histograms for
 742 the 16 corresponding SIFT subcells: for each $p \in \{1, \dots, 16\}$, to the corresponding
 743 p -th subcell s_j we associate the discrete histogram

$$744 \quad (54) \quad \tilde{H}_{j,\ell} = \frac{f_{16(p-1)+\ell}}{\sum_{\ell'=1}^8 f_{16(p-1)+\ell'}} \quad (1 \leq \ell \leq 8).$$

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

747 On Fig. 11, we display several reconstruction results obtained with the model
 748 MS-Poisson based on the multiscale HOGs or the true SIFT descriptors. As could

749 be expected, the reconstruction results obtained with the true SIFT descriptors are
 750 not as good as the ones obtained from multiscale HOGs, in particular many fine scale
 751 structures are lost, and the shape of small objects is not recovered in a coherent way
 752 (see for example the wings in the butterfly image). However, large-scale structures
 753 of the image are still retrieved quite properly which often suffices to understand the
 754 semantic content of the image.

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

764 **5.3. Comparison with previous works.** In this paragraph, we propose to
 765 compare our reconstruction models with the ones obtained by the methods by Wein-
 766 zaepfel et al. [64] and Dosovitskiy & Brox [19]. One important difference between
 767 these two other approaches and ours is that our method relies only on the content
 768 provided in the SIFT subcells while these methods exploit an external database ei-
 769 ther to copy local information from patches with similar SIFT descriptors (as in [64])
 770 or to build an up-convolutional neural network for reconstruction (as in [19]). Thus
 771 our work has no intention to outperform these methods in terms of visual quality of
 772 reconstruction (in particular, our method has absolutely no possibility of recovering
 773 the color information). Notice that we cannot compare to the method of [36] which is
 774 adapted to “dense SIFT” (i.e. SIFT descriptors computed on a dense set of patches)
 775 and not “sparse SIFT” (i.e. SIFT descriptors computed around the keypoints).

776 They are also minor differences in the extracted information because both these
 777 works do not rely on the original implementation of the SIFT method. The method
 778 of [64] actually uses “elliptic” interest regions (extracted using the Hessian-affine
 779 method by [42]) in which normalized multiscale HOG are computed (in the same
 780 way as in the original SIFT method). In contrast, Dosovitskiy and Brox use circular
 781 keypoints and descriptors that are computed with the VLFeat library [61]. But in
 782 order to apply an up-convolutional neural network to these features, they need to
 783 derive a grid-based representation of these features: the image is divided in 4×4
 784 cells and each cell containing a keypoint is being associated with the corresponding
 785 oriented keypoint and feature vector. If there is no keypoint, then they associate the
 786 zero vector, and if there are several keypoints they randomly choose one of them (see
 787 the details in [19, Section III]).

788 One advantage of the MS-Poisson model, compared to the result of [64], is that it
 789 is defined through the minimization of the global MS-Poisson energy (37). Therefore,
 790 it produces images that are globally coherent while respecting as much as possible the
 791 local constraints given by the multiscale HOGs. In contrast, the result of [64] is clearly
 792 affected by stitching artifacts which are inherent to their reconstruction method. On
 793 the other hand, their method is able to copy pieces of clean patches so that their
 794 reconstruction looks locally sharper (but also noisier).

795 However, the reconstructed images obtained in [19] are both globally coherent and
 796 quite sharp. Indeed, our method does not rely on an external database so it cannot
 797 compete with the one of [19], and in particular it cannot get back information which

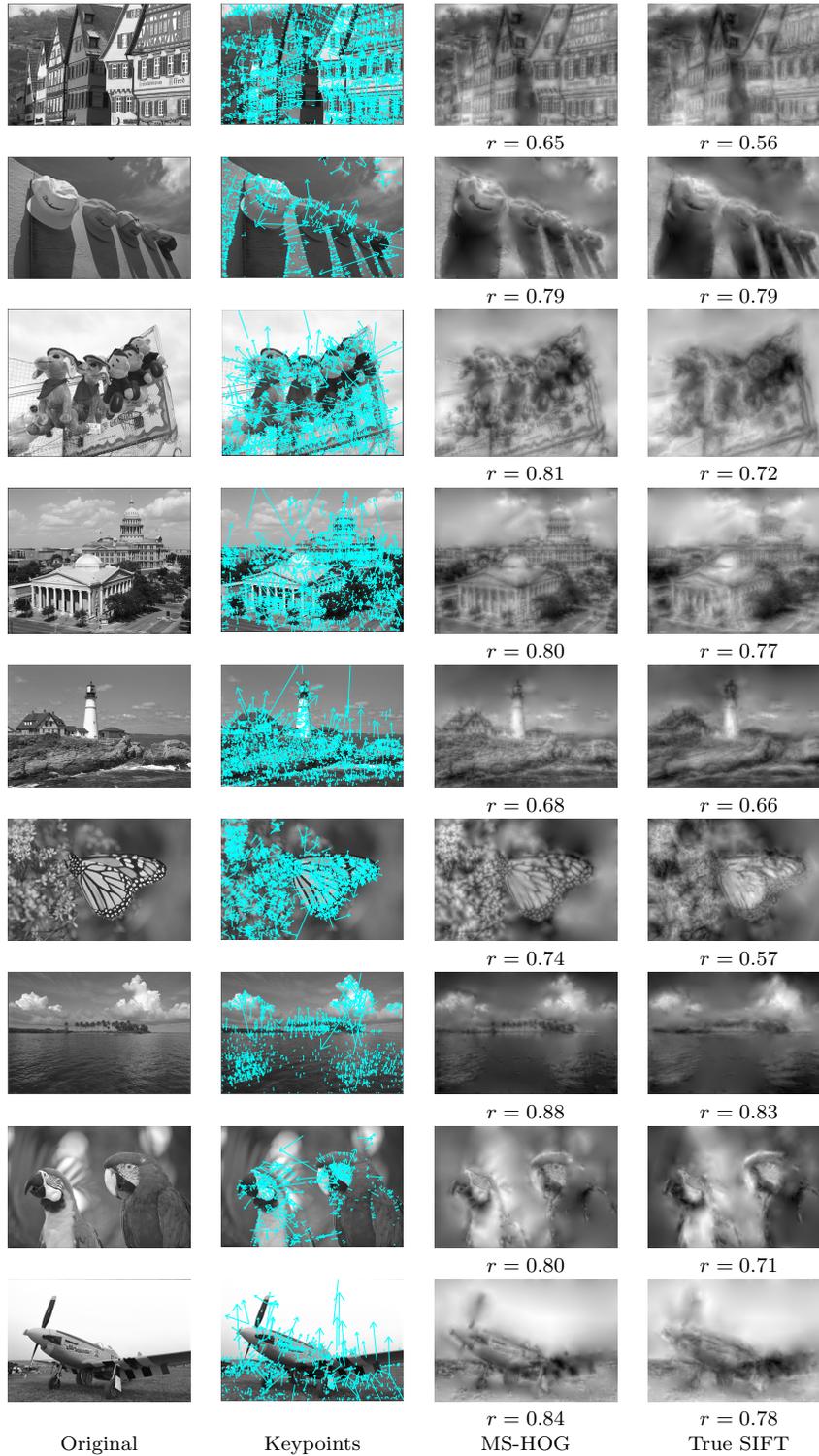


FIG. 11. *Reconstruction results from multiscale HOG or SIFT descriptors with images of the Live database [57]. For each row, from left to right, we display an original image, the same image with overimposed SIFT keypoints, a sample of the MS-Poisson model obtained from multiscale HOG, and a sample of 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.*

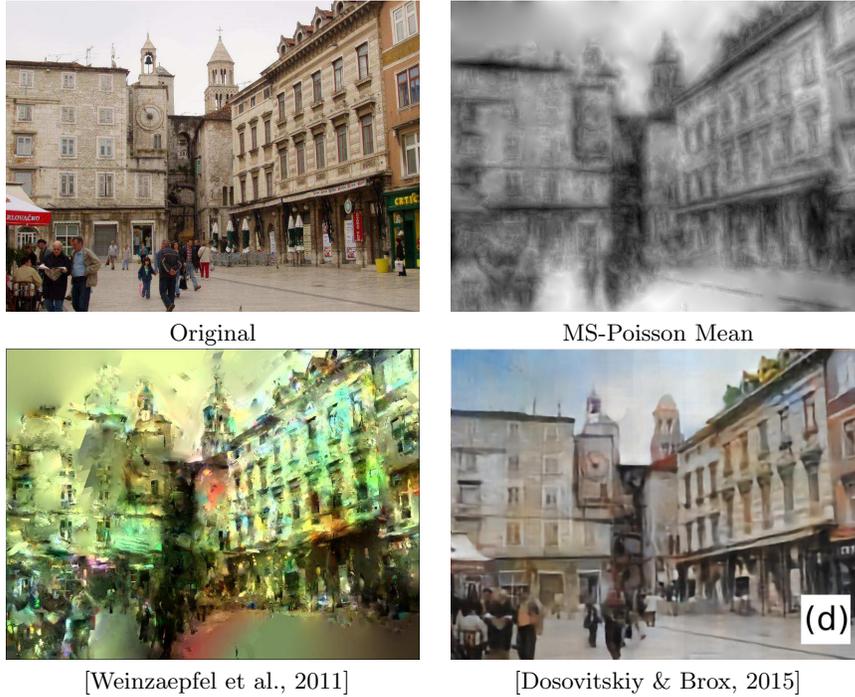


FIG. 12. *Comparison for SIFT reconstruction.* In the first row we display the original image and the reconstruction results obtained as the expectation of the MS-Poisson model computed on the true SIFT descriptors (see Section 5.2). In the second row we display the results obtained with the methods 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.

798 are completely lost in the SIFT descriptors (global contrast, or also color information).

799 **5.4. Reconstruction with other keypoints.** In this paragraph we question
 800 the very definition of the SIFT keypoints in terms of synthesis, in a similar way
 801 that what was done in [50]. Indeed, one can wonder if selecting the local extrema
 802 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
 803 extract relevant information for synthesis.

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

$$807 \quad (55) \quad \int u(\mathbf{x} + \mathbf{z}) g_\sigma(\mathbf{z}) d\mathbf{z} - u(\mathbf{x}) = \sigma^2 \Delta u(\mathbf{x}) + o(\sigma^2).$$

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

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

$$814 \quad (56) \quad (\mathbf{x}, \sigma) \mapsto |g_\sigma * u(\mathbf{x}) - u(\mathbf{x})|.$$

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

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

832 For these new sets of keypoints, we computed the average image of the MS-
 833 Poisson model. The results are displayed on Fig. 13. They clearly indicate that the
 834 usual SIFT keypoints lead to a reconstruction that is visually better than the oth-
 835 ers. The main problem of the “Min-Rec-Error” keypoints is that they do not extract
 836 enough small scale information: for the examples shown in Fig. 13 the average scale
 837 of these keypoints is approximately twice larger than the one of the SIFT keypoints.
 838 Besides, for both “Min-Rec-Error” and random keypoints, the spatial locations are
 839 not concentrated around geometric details as can be the case with the SIFT key-
 840 points. The comparison with “Random-grad” is particularly interesting: indeed the
 841 reconstruction with “Random-grad” keypoints is slightly better than the one with
 842 “Random-unif” keypoints, but still it fails to recover fine details. The main problem
 843 of the “Random-grad” approach is that it is not contrast invariant and thus it favors
 844 points with strong gradients in uniform regions over points in salient regions with low
 845 contrast. Thus, the usual definition of SIFT keypoints (and in particular the thresh-
 846 olding steps) is confirmed to be a relevant choice for extracting visual information
 847 near salient structures, both from the analysis or the synthesis perspective.

848 **6. Conclusion.** In this paper we proposed two stochastic models (MaxEnt, res-
 849 pectively MS-Poisson) for reconstructing an image based only on the information
 850 contained in the (monoscale, respectively multiscale) local HOGs computed in the
 851 SIFT subcells. With both models we get back images which are close to the original
 852 in terms of semantic content. This is still true if we compute the reconstructions based
 853 on the true SIFT descriptors. One benefit of these models over competing approaches
 854 is that they do not rely on any external image database, and besides the convolu-
 855 tive expressions found in this paper allow to compute statistics of the corresponding
 856 output random fields (e.g. local variance).

857 However, several questions raised by this work remain open. First it would be
 858 interesting to consider generalizations of the MS-Poisson model with different image
 859 priors, i.e. adopt other regularization terms in the functional. It is likely that solv-
 860 ing the corresponding optimization problem may require an iterative procedure, but
 861 on the other hand the solutions may exhibit cleaner geometric structures which are
 862 better extrapolated outside the SIFT subcells. Also, there is more to discuss about
 863 the optimality of keypoints with respect to the quality of reconstructed images. In

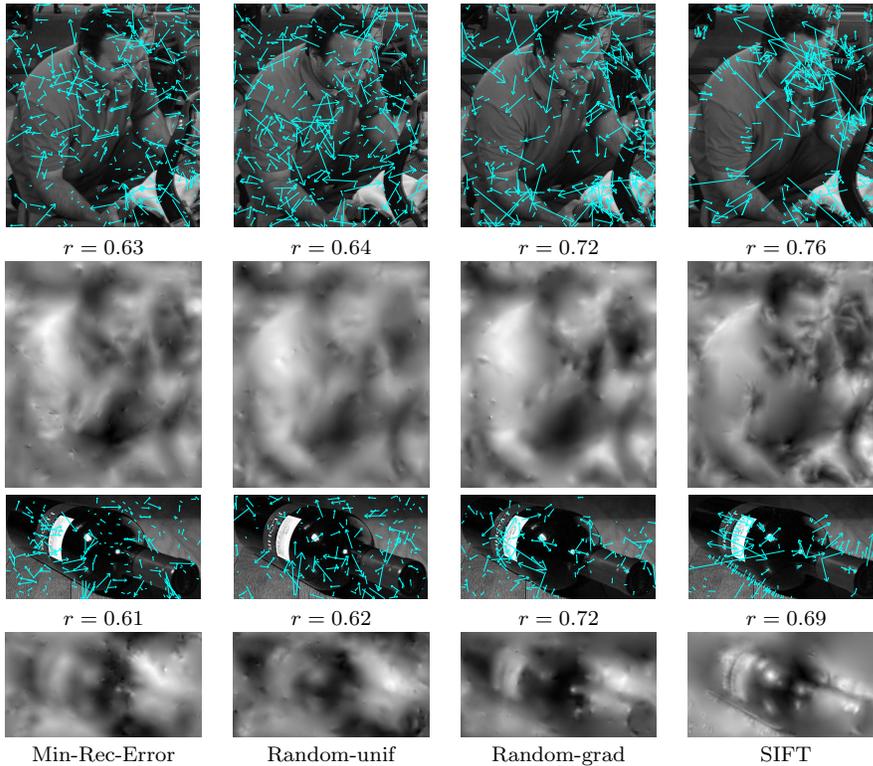


FIG. 13. **Reconstruction with other keypoints.** The first column (“Minimum reconstruction error”) corresponds to the keypoints obtained as local minima of (56). The second (“Random-unif”) and third column (“Random-grad”) corresponds to the randomly selected keypoints. The last column corresponds to the standard SIFT keypoints. The original images are displayed on Fig. 3 and Fig. 4. 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.

864 particular, here we adopted one unique reconstruction strategy in order to compare
 865 different sets of keypoints. But it seems possible to optimize both the sets of key-
 866 points and the reconstruction strategy in order to maximize a criterion linked to the
 867 proximity of the reconstruction to the input original image. This could be thought of
 868 as a kind of auto-encoding procedure in which the encoder is constrained to have a
 869 very particular form (that is, keypoint extractor).

870

REFERENCES

- 871 [1] T. AHONEN, A. HADID, AND M. PIETIKAINEN, *Face description with local binary patterns:*
 872 *Application to face recognition*, IEEE transactions on pattern analysis and machine intel-
 873 *ligence*, 28 (2006), pp. 2037–2041.
- 874 [2] A. ALAHI, R. ORTIZ, AND P. VANDERGHEYNST, *Freak: Fast retina keypoint*, in IEEE Conference
 875 *on Computer vision and pattern recognition (CVPR)*, IEEE, 2012, pp. 510–517.
- 876 [3] B. ALLEN AND M. KON, *The Marr Conjecture and Uniqueness of Wavelet Transforms*, arXiv
 877 *preprint arXiv:1401.0542*, (2015).
- 878 [4] B. ALLEN AND M. KON, *Unique recovery from edge information*, in *Sampling Theory and*
 879 *Applications (SampTA)*, 2015 International Conference on, IEEE, 2015, pp. 312–316.
- 880 [5] F. ATTNEAVE, *Some informational aspects of visual perception.*, *Psychological review*, 61
 881 (1954), p. 183.

- 882 [6] S. BATTIATO, G. GALLO, G. PUGLISI, AND S. SCCELLATO, *SIFT features tracking for video*
883 *stabilization*, in Image Analysis and Processing, 2007. ICIAP 2007. 14th International
884 Conference on, IEEE, 2007, pp. 825–830.
- 885 [7] H. BAY, A. ESS, T. TUYTELAARS, AND L. VAN GOOL, *Speeded-up robust features (SURF)*,
886 Computer vision and image understanding, 110 (2008), pp. 346–359.
- 887 [8] M. BLACK AND A. JEPSON, *Eigentracking: Robust matching and tracking of articulated objects*
888 *using a view-based representation*, International Journal of Computer Vision, 26 (1998),
889 pp. 63–84.
- 890 [9] Y.-L. BOUREAU, F. BACH, Y. LECUN, AND J. PONCE, *Learning mid-level features for recogni-*
891 *tion*, in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE,
892 2010, pp. 2559–2566.
- 893 [10] Y.-L. BOUREAU, N. LE ROUX, F. BACH, J. PONCE, AND Y. LECUN, *Ask the locals: multi-way*
894 *local pooling for image recognition*, in Computer Vision (ICCV), 2011 IEEE International
895 Conference on, IEEE, 2011, pp. 2651–2658.
- 896 [11] G. CSURKA, C. DANCE, L. FAN, J. WILLAMOWSKI, AND C. BRAY, *Visual categorization with*
897 *bags of keypoints*, in Workshop on statistical learning in computer vision, ECCV, 2004.
- 898 [12] S. CURTIS AND A. OPPENHEIM, *Reconstruction of multidimensional signals from zero crossings*,
899 J. Opt. Soc. Am. A, 4 (1987), pp. 221–231, doi:10.1364/JOSAA.4.000221.
- 900 [13] S. CURTIS, S. SHITZ, AND A. OPPENHEIM, *Reconstruction of nonperiodic two-dimensional sig-*
901 *nals from zero crossings*, IEEE Transactions on Acoustics, Speech, and Signal Processing,
902 35 (1987), pp. 890–893.
- 903 [14] N. DALAL AND B. TRIGGS, *Histograms of oriented gradients for human detection*, in Proceedings
904 of the IEEE CVPR, vol. 1, 2005, pp. 886–893.
- 905 [15] E. D'ANGELO, L. JACQUES, A. ALAHI, AND P. VANDERGHEYNST, *From bits to images: Inversion*
906 *of local binary descriptors*, IEEE Transactions on PAMI, 36 (2014), pp. 874–887.
- 907 [16] A. DESOLNEUX, *When the a contrario approach becomes generative*, International Journal of
908 Computer Vision, 116 (2016), pp. 46–65.
- 909 [17] A. DESOLNEUX AND A. LECLAIRE, *Stochastic image reconstruction from local histograms of gra-*
910 *dient orientation*, in Proceedings of the sixth International Conference on Scale Space and
911 Variational Methods in Computer Vision (SSVM), Springer, Lecture Notes in Computer
912 Science, 2017, pp. 133–145.
- 913 [18] A. DESOLNEUX, L. MOISAN, AND J. MOREL, *From Gestalt theory to image analysis: a proba-*
914 *bilistic approach*, vol. 34, Springer Science & Business Media, 2007.
- 915 [19] A. DOSOVITSKIY AND T. BROX, *Inverting Visual Representations with Convolutional Networks*,
916 arXiv:1506.02753 [cs], (2015).
- 917 [20] J. H. ELDER AND S. W. ZUCKER, *Scale space localization, blur, and contour-based image coding*,
918 in Computer Vision and Pattern Recognition, 1996. Proceedings CVPR'96, 1996 IEEE
919 Computer Society Conference on, IEEE, 1996, pp. 27–34.
- 920 [21] O. FAUGERAS, *Three-dimensional computer vision: a geometric viewpoint*, MIT press, 1993.
- 921 [22] P. FELZENSZWALB, R. GIRSHICK, D. MCALLESTER, AND D. RAMANAN, *Object detection with dis-*
922 *criminatively trained part-based models*, IEEE Transactions on PAMI, 32 (2010), pp. 1627–
923 1645.
- 924 [23] C. HARRIS AND M. STEPHENS, *A combined corner and edge detector.*, in Alvey vision conference,
925 vol. 15, Citeseer, 1988, p. 50.
- 926 [24] R. HUMMEL AND R. MONIOT, *Reconstructions from zero crossings in scale space*, IEEE Trans-
927 actions on Acoustics, Speech, and Signal Processing, 37 (1989), pp. 2111–2130.
- 928 [25] F. JUEFEI-XU AND M. SAVVIDES, *Learning to invert local binary patterns.*, in BMVC, 2016.
- 929 [26] S. K. AND Z. A., *Very deep convolutional networks for large-scale image recognition*, in Pro-
930 ceedings of the International Conference on Learning Representations, 2014.
- 931 [27] H. KATO AND T. HARADA, *Image reconstruction from bag-of-visual-words*, in Proceedings of
932 the IEEE CVPR, 2014, pp. 955–962.
- 933 [28] A. KRIZHEVSKY, I. SUTSKEVER, AND G. HINTON, *Imagenet classification with deep convolutional*
934 *neural networks*, in Advances in neural information processing systems, 2012, pp. 1097–
935 1105.
- 936 [29] S. LAZEBNIK, C. SCHMID, AND J. PONCE, *Beyond bags of features: Spatial pyramid matching*
937 *for recognizing natural scene categories*, in IEEE computer society conference on Computer
938 vision and pattern recognition, vol. 2, IEEE, 2006, pp. 2169–2178.
- 939 [30] S. LEUTENEGGER, M. CHLI, AND R. Y. SIEGWART, *Brisk: Binary robust invariant scalable*
940 *keypoints*, in IEEE International Conference on Computer Vision (ICCV), IEEE, 2011,
941 pp. 2548–2555.
- 942 [31] T. LINDBERG, *Feature detection with automatic scale selection*, International journal of com-
943 puter vision, 30 (1998), pp. 79–116.

- 944 [32] T. LINDBERG, *Image matching using generalized scale-space interest points*, Journal of Math-
945 ematical Imaging and Vision, 52 (2015), pp. 3–36.
- 946 [33] D. LOWE, *Distinctive image features from scale-invariant keypoints*, International Journal of
947 Computer Vision, 60 (2004), pp. 91–110.
- 948 [34] Y. LU, S. ZHU, AND Y. N. WU, *Learning FRAME models using CNN filters for*
949 *knowledge visualization*, CoRR, abs/1509.08379 (2015), <http://arxiv.org/abs/1509.08379>,
950 [arXiv:1509.08379](https://arxiv.org/abs/1509.08379).
- 951 [35] A. MAHENDRAN AND A. VEDALDI, *Understanding deep image representations by inverting them*,
952 in IEEE CVPR, 2015, pp. 5188–5196.
- 953 [36] A. MAHENDRAN AND A. VEDALDI, *Visualizing deep convolutional neural networks using natural*
954 *pre-images*, International Journal of Computer Vision, 120 (2016), pp. 233–255.
- 955 [37] E. MAIR, G. D. HAGER, D. BURSCHKA, M. SUPPA, AND G. HIRZINGER, *Adaptive and generic*
956 *corner detection based on the accelerated segment test*, in European conference on Com-
957 puter vision, Springer, 2010, pp. 183–196.
- 958 [38] S. MALLAT AND S. ZHONG, *Characterization of signals from multiscale edges*, IEEE Trans-
959 actions on PAMI, 14 (1992), pp. 710–732.
- 960 [39] D. MARR, *Vision: A computational investigation into the human representation and processing*
961 *of visual information*, W.H. Freeman and Company, 1982.
- 962 [40] D. MARR AND E. HILDRETH, *Theory of edge detection*, Proceedings of the Royal Society of
963 London B: Biological Sciences, 207 (1980), pp. 187–217.
- 964 [41] Y. MEYER, *Wavelets-algorithms and applications*, vol. 1, Society for Industrial and Applied
965 Mathematics Translation, 1993.
- 966 [42] K. MIKOLAJCZYK AND C. SCHMID, *Scale & affine invariant interest point detectors*, Interna-
967 tional journal of computer vision, 60 (2004), pp. 63–86.
- 968 [43] K. MIKOLAJCZYK AND C. SCHMID, *A performance evaluation of local descriptors*, IEEE Trans-
969 actions on PAMI, 27 (2005), pp. 1615–1630.
- 970 [44] J.-M. MOREL, A. PETRO, AND C. SBERT, *Fourier implementation of Poisson image editing*,
971 Pattern Recognition Letters, 33 (2012), pp. 342–348.
- 972 [45] J.-M. MOREL AND G. YU, *ASIFT: A new framework for fully affine invariant image compar-*
973 *ison*, SIAM Journal on Imaging Sciences, 2 (2009), pp. 438–469.
- 974 [46] J.-M. MOREL AND G. YU, *Is SIFT scale invariant?*, Inverse Problems and Imaging, 5 (2011),
975 pp. 115–136.
- 976 [47] D. MUMFORD AND A. DESOLNEUX, *Pattern Theory: The Stochastic Analysis of Real-World*
977 *Signals*, A K Peters/CRC Press, Natick, Mass, 2010.
- 978 [48] P. MUSÉ, F. SUR, F. CAO, Y. GOUSSEAU, AND J.-M. MOREL, *An a contrario decision method for*
979 *shape element recognition*, International Journal of Computer Vision, 69 (2006), pp. 295–
980 315.
- 981 [49] Y. NESTEROV, *Introductory lectures on convex optimization: A basic course*, vol. 87, Springer,
982 2004.
- 983 [50] M. NIELSEN AND M. LILLHOLM, *What do features tell about images?*, in Scale-Space, vol. 1,
984 Springer, 2001, pp. 39–50.
- 985 [51] T. OJALA, M. PIETIKÄINEN, AND T. MÄENPÄÄ, *Multiresolution gray-scale and rotation in-*
986 *variant texture classification with local binary patterns*, IEEE Transactions on PAMI, 24
987 (2002), pp. 971–987.
- 988 [52] P. PÉREZ, M. GANGNET, AND A. BLAKE, *Poisson Image Editing*, in ACM SIGGRAPH 2003
989 Papers, SIGGRAPH '03, 2003, pp. 313–318, [doi:10.1145/1201775.882269](https://doi.org/10.1145/1201775.882269).
- 990 [53] J. PHILBIN, O. CHUM, M. ISARD, J. SIVIC, AND A. ZISSERMAN, *Object retrieval with large vo-*
991 *cabularies and fast spatial matching*, in IEEE Conference on Computer Vision and Pattern
992 Recognition, 2007., IEEE, 2007, pp. 1–8.
- 993 [54] I. REY OTERO AND M. DELBRACIO, *Anatomy of the SIFT Method*, Image Processing On Line,
994 4 (2014), pp. 370–396, [doi:10.5201/ipol.2014.82](https://doi.org/10.5201/ipol.2014.82).
- 995 [55] E. ROSTEN AND T. DRUMMOND, *Machine learning for high-speed corner detection*, in Computer
996 Vision—ECCV 2006, Springer, 2006, pp. 430–443.
- 997 [56] J. SANZ AND T. HUANG, *Theorems and experiments on image reconstruction from zero cross-*
998 *ings*, 1987. IBM Almaden Research Center.
- 999 [57] H. SHEIKH, Z. WANG, L. CORMACK, AND A. BOVIK, *Live image quality assessment database*
1000 *release 2 (2005)*, 2005.
- 1001 [58] J. SIVIC AND A. ZISSERMAN, *Video Google: A text retrieval approach to object matching in*
1002 *videos*, in Proceedings of the IEEE ICCV, 2003, pp. 1470–1477.
- 1003 [59] T. TUYTELAARS AND K. MIKOLAJCZYK, *Local invariant feature detectors: a survey*, Foundations
1004 and trends in computer graphics and vision, 3 (2008), pp. 177–280.
- 1005 [60] T. TUYTELAARS AND L. VAN GOOL, *Matching widely separated views based on affine invariant*

- 1006 *regions*, International journal of computer vision, 59 (2004), pp. 61–85.
- 1007 [61] A. VEDALDI AND B. FULKERSON, *Vlfeat: An open and portable library of computer vision*
1008 *algorithms*, in Proceedings of the 18th ACM international conference on Multimedia, ACM,
1009 2010, pp. 1469–1472.
- 1010 [62] C. VONDRICK, A. KHOSLA, T. MALISIEWICZ, AND A. TORRALBA, *Hoggles: Visualizing object*
1011 *detection features*, in Proceedings of the IEEE ICCV, 2013, pp. 1–8.
- 1012 [63] C. WALLRAVEN, B. CAPUTO, AND A. GRAF, *Recognition with local features: the kernel recipe*,
1013 in Proceedings of the IEEE ICCV, 2003, pp. 257–264.
- 1014 [64] P. WEINZAEPFEL, H. JÉGOU, AND P. PÉREZ, *Reconstructing an image from its local descriptors*,
1015 in Proceedings of the IEEE CVPR, 2011, pp. 337–344.
- 1016 [65] J. YANG, D. SCHONFELD, AND M. MOHAMED, *Robust video stabilization based on particle filter*
1017 *tracking of projected camera motion*, IEEE Transactions on Circuits and Systems for Video
1018 Technology, 19 (2009), pp. 945–954.
- 1019 [66] A. YILMAZ, O. JAVED, AND M. SHAH, *Object tracking: A survey*, ACM computing surveys
1020 (CSUR), 38 (2006), p. 13.
- 1021 [67] M. ZEILER AND R. FERGUS, *Visualizing and understanding convolutional networks*, in Proceed-
1022 *ings of ECCV*, Springer, 2014, pp. 818–833.
- 1023 [68] J. ZHANG, M. MARSZALEK, S. LAZEBNIK, AND C. SCHMID, *Local features and kernels for clas-*
1024 *sification of texture and object categories: A comprehensive study*, International journal of
1025 computer vision, 73 (2007), pp. 213–238.
- 1026 [69] S. ZHU, Y. WU, AND D. MUMFORD, *Filters, random fields and maximum entropy (FRAME):*
1027 *Towards a unified theory for texture modeling*, International Journal of Computer Vision,
1028 27 (1998), pp. 107–126.