



# A multi-scale probabilistic network model for detection, synthesis and compression in mammographic image analysis

Paul Sajda<sup>a,\*</sup>, Clay Spence<sup>b</sup>, Lucas Parra<sup>b</sup>

<sup>a</sup>Department of Biomedical Engineering, Columbia University, 351 Engineering Terrace, MC 8904, New York, NY 10027, USA

<sup>b</sup>Sarnoff Corporation, 201 Washington Road, Princeton, NJ, USA

Accepted 7 January 2003

## Abstract

We develop a probabilistic network model over image spaces and demonstrate its broad utility in mammographic image analysis, particularly with respect to computer-aided diagnosis. The model employs a multi-scale pyramid decomposition to factor images across scale and a network of tree-structured hidden variables to capture long-range spatial dependencies. This factoring makes the computation of the density functions local and tractable. The result is a hierarchical mixture of conditional probabilities, similar to a hidden Markov model on a tree. The model parameters are found with maximum likelihood estimation using the expectation-maximization algorithm. The utility of the model is demonstrated for three applications: (1) detection of mammographic masses for computer-aided diagnosis; (2) qualitative assessment of model structure through mammographic synthesis; and (3) compression of mammographic regions of interest. © 2003 Elsevier Science B.V. All rights reserved.

**Keywords:** Probabilistic network model; Multi-scale pyramid decomposition; Mammographic computer-aided diagnosis; Image synthesis; Image compression

## 1. Introduction

Computer-aided diagnosis (CAD) can be defined as a diagnosis made by a radiologist who incorporates the results of computer analysis of the radiographs (Doi et al., 1993). The goal of CAD is to improve radiologists' performance by indicating the sites of potential abnormalities, to reduce the number of missed lesions, and/or by providing quantitative analysis of specific regions in an image to improve diagnosis. CAD systems typically operate as automated "second-opinion" or "double-reading" systems that indicate lesion location and/or type. Since individual human observers overlook different findings, it has been shown that double reading (the review of a study by more than one observer) increases the detection rate of breast cancers by 5–15% (Bird, 1990; Metz and Shen,

1992; Thurfjell et al., 1994). Double reading, if not done efficiently, can significantly increase the cost of screening, given the need for a second radiologist/mammographer. Methods to provide improved detection with little increase in cost will have significant impact on the benefits of screening. Automated CAD systems are a promising approach for low-cost double-reading.

Several CAD systems have been developed for mammographic screening and the first have been approved by the FDA.<sup>1</sup> Complete systems have been rigorously characterized, both in retrospective and prospective trials (Burhenne et al., 2000). Though many have demonstrated clinical utility, there is still a need to reduce false-positive rates generated by CAD systems. For example, prospective clinical studies have shown lower sensitivities and specificities than originally found in retrospective studies—80%

\*Corresponding author. Tel.: +1-212-854-5270.  
E-mail address: [sajda@columbia.edu](mailto:sajda@columbia.edu) (P. Sajda).

<sup>1</sup>R2 Technology M1000 Image Checker, CADx Medical Second-Look, and Intelligent Systems Software MammoReader.

62 cancers detected with 2.4 false positives per case in  
 63 prospective studies versus 85–90% sensitivity at one to  
 64 two false positives per case in retrospective studies  
 65 (Nishikawa et al., 1996).

### 66 1.1. The role of statistical pattern recognition and 67 neural networks in CAD

68 CAD systems usually consist of two distinct subsystems,  
 69 one designed to detect microcalcifications and one to  
 70 directly detect masses (Giger et al., 2000). A common  
 71 element in both subsystems is a statistical pattern recogni-  
 72 tion model, used to improve detection and reduce false-  
 73 positive rates introduced by earlier stages of processing.  
 74 Neural networks are a particularly important class of  
 75 statistical model in CAD because they are able to capture  
 76 complicated, often nonlinear, relationships in high dimen-  
 77 sional feature spaces not easily captured by heuristic or  
 78 rule-based algorithms. Several groups have developed  
 79 neural networks architectures for CAD. Many of these  
 80 architectures exploit well-known features that might also  
 81 be used by radiologists (Floyd et al., 1994; Jiang et al.,  
 82 1996; Huo et al., 1998), while others utilize more generic  
 83 feature sets (Zhang et al., 1994; Lo et al., 1996; Chan et  
 84 al., 1998; Sajda et al., 2002). In general, these neural  
 85 networks are *recognition* or *discriminative* probabilistic  
 86 models (Dayan and Abbott, 2002) in that they estimate  
 87  $\Pr(C | I)$ , the conditional probability of class  $C$  (e.g., mass  
 88 vs. non-mass) given image  $I$  or a set of features extracted  
 89 from  $I$ . An alternative approach is to construct a *generative*  
 90 probabilistic model of the data, which, using the afore-  
 91 mentioned formulation, would be a model that estimates  
 92 the class conditional distribution,  $\Pr(I | C)$ . Such a model  
 93 has several attractive features for mammographic image  
 94 analysis. For example, classification is possible by training  
 95 a distribution for each class and using Bayes' rule to obtain  
 96  $\Pr(C | I) = \Pr(I | C)\Pr(C)/\Pr(I)$ . In addition, novel exam-  
 97 ples, relative to the training data used to build the model,  
 98 can be detected by computing the absolute likelihood over  
 99 each model. In terms of CAD, the ability to identify novel  
 100 examples is useful for establishing confidence measures on  
 101 the CAD output (e.g., should the output of the classifier be  
 102 "trusted" given that the current data is very different from  
 103 the training data). In addition, novelty detection can be  
 104 used to identify new clinical data that might be used to  
 105 re-train/refine the CAD system. Since essentially any type  
 106 of image analysis can be formulated given knowledge of  
 107 the distribution of the data, the generative probabilistic  
 108 model can also be used to compress (Cover and Thomas,  
 109 1991), suppress noise (Romberg et al., 2001), interpolate,  
 110 increase or extend resolution (Freeman et al., 2002), etc.

### 111 1.2. Generative probabilistic models for images

112 Previous research has focused on developing probabilis-  
 113 tic models of biological and natural shapes, much of which

is based on the work of Grenander et al. (1991). This has  
 in turn led to the development of active shape and  
 appearance models for medical image analysis, most  
 notably those of Cootes and Taylor (Cootes et al., 1994;  
 Cootes and Taylor, 2001). These approaches construct a  
 statistical description of object shape over a set of land-  
 marks that are often extracted by a human expert (e.g.,  
 radiologist). The statistical descriptions are formulated as  
 generative models and can be sampled to construct new  
 instances of a given shape. The approach has demonstrated  
 great utility in localizing structure in medical imagery,  
 particularly in cases where the structure is well described  
 by its contours/borders (e.g., ventricles in the brain and  
 heart). However, in the case of mammography, lesions are  
 not well characterized by border shape alone. Radiologists  
 typically integrate evidence which includes texture, homo-  
 geneity, and spiculation, as well as contextual information  
 such as vascularization and proximity to mammillary ducts  
 (Kopans, 1989). Therefore, a concise set of landmarks is  
 not easily extracted and instead a classification system  
 must learn the set of shape and non-shape features which  
 are correlated with disease (or absence of disease).

Significant efforts have also focused on the construction  
 of generative probabilistic models for directly modeling  
 images. Grenander (1983) was one of the first to propose a  
 Bayesian framework for image analysis. This framework  
 led to the development of a series of image distribution  
 models, most notably the Markov Random Field (MRF)  
 developed by Geman and Geman (1984) and further  
 developed and studied by others (e.g., Chellappa and  
 Chatterjee, 1985). MRFs model distributions by assuming  
 that images are locally smooth except for relatively sparse  
 intensity gradients and edges. The underlying assumption  
 of an MRF is that local image structure is sufficient for  
 global image representation. However, these models tend  
 to be computationally expensive, have limited forms for  
 the distributions/potential functions, and have difficulty  
 capturing more global structure and long-range dependen-  
 cies in images. Zhu et al. (1997) attempted to overcome  
 some of the limitations in MRFs by computing distribu-  
 tions over a set of features constructed from the histograms  
 of filtered images (e.g., using Gabor filters). In their  
 approach they compute the maximum entropy distribution  
 given the statistics across these features. Though this  
 approach works well for textures, it is not clear how well it  
 models the appearance of more structured objects.

De Bonet and Viola proposed a flexible histogram  
 approach (De Bonet and Viola, 1998; De Bonet et al.,  
 1998), where features are extracted at multiple image  
 scales, with the resulting feature vectors treated as a set of  
 independent samples drawn from a distribution. The  
 distribution of feature vectors is subsequently modeled  
 using Parzen windows. Though they report good results,  
 their model treats the feature vectors from neighboring  
 pixels as independent samples, when in fact they share  
 exactly the same components from lower resolutions. One

173 solution to this is to build a model in which the features at  
 174 one pixel of one pyramid level condition the features at  
 175 each of several child pixels at the next higher-resolution  
 176 pyramid level. The multiscale stochastic process (MSP)  
 177 methods do exactly that. Luetngen and Willsky (1995), for  
 178 example, applied a scale-space auto-regression (AR) model  
 179 to texture discrimination. They use a quadtree or quadtree-  
 180 like organization of the pixels in an image pyramid, and  
 181 model the features in the pyramid as a stochastic process  
 182 from coarse-to-fine levels along the tree. The variables in  
 183 the process are hidden, and the observations are sums of  
 184 these hidden variables plus noise. However, the assumed  
 185 Gaussian distributions are a limitation of MSP models as  
 186 well as the fact that the model is of the probability of the  
 187 observations on the tree, not of the image. Once again,  
 188 these methods appear well suited for modeling texture, but  
 189 it is unclear how one might build models to capture the  
 190 appearance of more structured objects, such as mammog-  
 191 raphic masses.

192 Recently, several groups have developed what are  
 193 essentially extensions of the MSP model by adding hidden  
 194 variables. These can be seen as improving the model's  
 195 ability to capture non-local dependencies in the image. For  
 196 example, Crouse et al. (1998) developed their Hidden  
 197 Markov Tree (HMT) models for signals and images. A  
 198 primary motivation of these models is to capture the  
 199 tendency for wavelet coefficients to group into two classes,  
 200 one with large and the other with small coefficient  
 201 magnitudes. Thus their hidden states have one of two  
 202 values corresponding to large and small wavelet coeffi-  
 203 cients. This is well suited to the many signal and image  
 204 types that have homogeneous regions with boundaries.  
 205 These models have been successfully applied to several  
 206 problems, especially image enhancement and texture seg-  
 207 mentation (Romberg et al., 2001; Coi and Baraniuk, 2001).  
 208 Cheng and Bouman (2001) applied a similar model for  
 209 segmentation, in which the observed class labels play the  
 210 role of hidden variables, and therefore are no longer  
 211 hidden.

212 We have developed a class of models for probability  
 213 distributions of images that we call hierarchical image  
 214 probability (HIP) models. The HIP model can be viewed  
 215 as a development of the HMT model, with several  
 216 differences. The main elements of both the HIP and HMT  
 217 models include:

- 219 • Capturing local dependencies in a coarse-to-fine factor-  
 220 ing of the image distribution over scale and position.
- 221 • Capturing non-local and scale dependencies through a  
 222 set of discrete hidden variables whose dependency  
 223 graph is a tree.
- 224 • Optimizing model parameters to match the natural  
 225 image statistics using strict Maximum Likelihood.
- 226 • Enabling both evaluation of the likelihood and sampling  
 227 from the distribution.

228 In addition, the HIP model differs from the HMT model  
 in the following ways:

- The coefficients of the different subbands at each node  
 are modeled jointly, using mixtures of multivariate  
 Gaussian distributions.<sup>2</sup>
- The number of hidden states in each level is adjusted  
 separately in an attempt to better fit the image dis-  
 tribution.
- The hidden states capture complex structure in the  
 image through the use of mixture, hierarchy and scale  
 components.
- The probability of a child state value at a node,  
 conditioned on the state at the parent node, also  
 depends on the child node's relative position, e.g.  
 upper-left, lower-right, etc.
- The mean of each normal distribution depends on the  
 corresponding coefficient vector in the unsampled  
 wavelet coefficient subbands from the next coarsest  
 resolution pyramid level. (The HIP model resembles a  
 simple MSP model in this way.)

247 In the following we begin by presenting the structure of  
 248 the HIP model, along with an EM algorithm used to  
 249 estimate its parameters. We first describe the most simple  
 250 form of the HIP model, namely with a single component in  
 251 the hidden variable structure, and then augment the model  
 252 to include mixture, hierarchy and scale components. We  
 253 then demonstrate the broad utility of the complete model  
 254 by presenting results for several applications in mammog-  
 255 raphic image analysis, including mass detection in CAD,  
 256 mammographic synthesis, and compression of mammog-  
 257 raphic ROIs. In all cases we compare results to a tradition-  
 258 al HMT (Crouse et al., 1998).

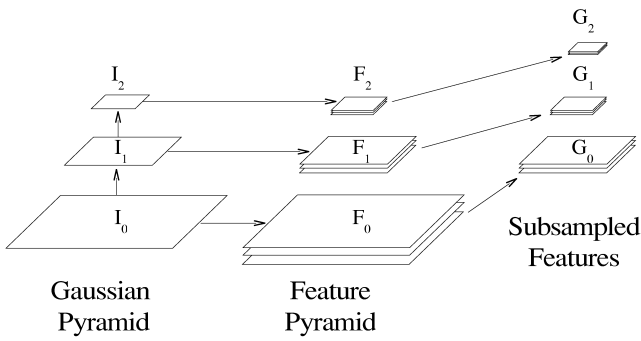
## 2. Structure of the HIP model

### 2.1. Coarse-to-fine factoring of image distributions

261 Similar to previous work (Luetngen and Willsky, 1995;  
 262 De Bonet and Viola, 1998; De Bonet et al., 1998; Crouse  
 263 et al., 1998; Cheng and Bouman, 2001) we model depen-  
 264 dencies in an image over a range of scales. We begin by  
 265 representing the image as a set of feature images, for  
 266 example computed using a set of filters with limited spatial  
 267 support. Coarse-scale image structure is captured by  
 268 applying the set of filters at a low-resolution in a pyramid  
 269 decomposition of the image. Long-range dependencies of  
 270 fine-scale structure are modeled by conditioning fine scales  
 271 on coarse scales. Denoting the set of feature images at  
 272 pyramid level  $l$  by  $\mathbf{F}_l$ , the goal is to write the image  
 273 distribution in a form similar to  $\Pr(I) \sim$   
 274  $\Pr(\mathbf{F}_0 | \mathbf{F}_1) \Pr(\mathbf{F}_1 | \mathbf{F}_2) \dots$

275 We first prove that a coarse-to-fine factoring of this form  
 276 is exact. From an image  $I$  build a low-pass (e.g., Gaussian)  
 277 pyramid. Call the  $l$ th level of this pyramid  $I_l$ , with the

<sup>2</sup>Arbitrarily complex distributions can be modeled as mixtures of  
 Gaussians.



280

281 Fig. 1. Pyramids and feature notation used to demonstrate coarse-to-fine  
282 factoring.

288 original full resolution image denoted as  $I_0$  (see Fig. 1).  
289 For each low-pass image  $I_l$  at level  $l$  extract a set of feature  
290 images  $F_l$ . Appropriate sub-sampling of these features  
291 results in  $G_l$ , having the same dimensions as  $I_{l+1}$ . Denote  
292 by  $\tilde{G}_l$  the set of images containing  $I_{l+1}$  and the images in  
293  $G_l$ . Finally, denote the mapping from  $I_l$  to  $\tilde{G}_l$  as  $\tilde{\mathcal{G}}_l$ . We  
294 continue this decomposition up to layer  $L$ .<sup>3</sup> Wavelet  
295 transforms are well-known examples of such mappings.  
296 Alternately, one could build Gaussian pyramids (Burt and  
297 Adelson, 1983) to obtain  $I_l$  and then filter these with  
298 several carefully chosen band-pass filters, followed by  
299 subsampling, as shown in Fig. 1.

300 Suppose now that  $\tilde{\mathcal{G}}_l : I_l \rightarrow \tilde{G}_l$  is invertible. Then  $\tilde{\mathcal{G}}_l$  is a  
301 change of variables, and we can relate the distributions on  
302 the two different sets of variables through multiplication  
303 by the Jacobian, i.e.,  $\Pr(I_l) = |\tilde{\mathcal{G}}_l| \Pr(\tilde{G}_l)$ . Since  $\tilde{G}_l =$   
304  $(G_l, I_{l+1})$ ,  $\Pr(\tilde{G}_l)$  can be factored to obtain  $\Pr(I_l) =$   
305  $|\tilde{\mathcal{G}}_l| \Pr(G_l | I_{l+1}) \Pr(I_{l+1})$ . If  $\tilde{\mathcal{G}}_l$  is invertible for all  $l \in$   
306  $\{0, \dots, L\}$  then we can recursively apply this change of  
307 variables and factoring procedure to obtain<sup>4</sup>

$$308 \Pr(I) = \left[ \prod_{l=0}^L |\tilde{\mathcal{G}}_l| \Pr(G_l | I_{l+1}) \right]. \quad (1)$$

309 This is a very general result, valid for all  $\Pr(I)$ , requiring  
310 only that the mapping be invertible and unique.

311 If our features are the outputs of linear filters, the  
312 determinants  $|\tilde{\mathcal{G}}_l|$  depend only on the filters used, and not  
313 on the image or model parameters. Therefore, we can drop  
314 the determinants if we write Eq. (1) as a proportionality

$$315 \Pr(I) \propto \prod_{l=0}^L \Pr(G_l | I_{l+1}). \quad (2)$$

316 Note that for comparing the likelihoods with different

283 <sup>3</sup>It will prove convenient to define  $G_L$  to be the same as  $\tilde{G}_L$ .

284 <sup>4</sup>For the last layer  $L$  the conditioning on  $I_{L+1}$  is to be ignored, since we  
285 defined  $G_L$  to include  $I_{L+1}$ .

features or sampling from the distribution (e.g., synthesis) 317  
it will be important to keep the determinants. 318

## 2.2. Hidden variables for modeling non-local dependencies 319

320 For the sake of computational tractability we would like  
321 to factor  $\Pr(G_l | I_{l+1})$  over position. However, this is  
322 problematic due to non-local dependencies that remain  
323 after coarse-to-fine factoring. Fig. 2 illustrates these dependen-  
324 cies. Assume that the presence of a particular object,  $O_A$   
325 (e.g., mammographic mass), may be inferred with high  
326 probability at a coarse scale,  $l + 1$ , of the image pyramid.  
327 Assume further that the presence of  $O_A$  implies with high  
328 probability the presence of a texture at position  $x$  at a finer  
329 scale  $l$ . Since we may need to examine an extended spatial  
330 area at level  $l + 1$  to detect the object, the presence of the  
331 texture at  $x$  in level  $l$  can depend upon an extended spatial  
332 area in level  $l + 1$ , i.e., the dependence between scales is  
333 non-local. Similarly (see Fig. 2(B)), the information at  
334 level  $l + 1$  may be sufficient for detecting an object but not  
335 for distinguishing its class,  $O_A$  or  $O_B$  (e.g., mass vs.  
336 non-mass). However, the detection of a single object  
337 imposes constraints of structure at high resolutions, for  
338 instance that distant positions in  $l$  have similar texture.  
339 Once again, conditioning fine scales on coarser scales  
340 cannot capture these long-range dependencies, which are  
341 entirely within the finer scale in this example. 342

343 To capture these dependencies we introduce a hidden  
344 variable  $A$  that takes on values in some set  $\mathcal{A}$ . We assume  
345 that  $A$  contains sufficient information for  $\Pr(G_l | I_{l+1}, A)$   
346 to factor over position  $x$

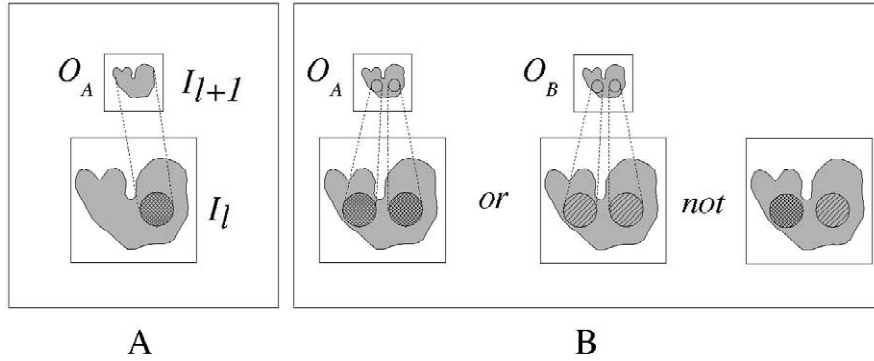
$$347 \Pr(I) \propto \sum_{A \in \mathcal{A}} \left[ \prod_{l=0}^L \Pr(G_l | I_{l+1}, A) \right] \Pr(A) \quad (3)$$

$$= \sum_{A \in \mathcal{A}} \left[ \prod_{l=0}^L \prod_{x \in \mathcal{P}_l} \Pr(\mathbf{g}_l(x) | I_{l+1}, A) \right] \Pr(A).$$

348 Here,  $\mathcal{P}_l$  is the set of all positions in resolution level  $l$   
349 in the wavelet/multi-resolution decomposition.<sup>5</sup> Note that by  
350 replacing uppercase letters (e.g.,  $G_l$ ) with lowercase letters  
351 which are functions of  $x$  (e.g.,  $\mathbf{g}_l(x)$ ) we are indicating a  
352 factoring of the features over position.

353 To simplify, we assume that, given  $A$ ,  $\mathbf{g}_l(x)$  depends  
354 only on the local information in  $I_{l+1}$  which is captured by  
355  $\mathbf{f}_{l+1}(x)$ , the features of  $I_{l+1}$  at position  $x$ . To be precise, the  
356 complete decomposition of  $I_{l+1}$  requires in addition to the  
357 high-pass features  $F_{l+1}$  also the low-pass information  $I_{l+2}$ .  
358 To simplify the presentation, we drop this, essentially  
359 assuming that  $A$  carries all of the coarse-scale intensity  
360 information from  $I_{l+2}$  that is relevant for  $G_l$ . (In practice, it

286 <sup>5</sup>In the following we frequently simplify expressions by omitting the  
287 limits of the sums and products, since they should be clear from context.



363

364 Fig. 2. Example of dependencies that cannot be captured by a coarse-to-fine factoring. (A) An object  $O_A$ , detectable at level  $l + 1$ , implies the presence of  
 365 particular texture at location  $x$  at level  $l$ . Since we place no constraints on the spatial extent of  $O_A$  in level  $l + 1$ , the presence of the texture at  $x$  can depend  
 366 upon an extended region in  $l + 1$ . (B) The information in level  $l + 1$  may be insufficient to discriminate between objects  $O_A$  and  $O_B$ , however the detection  
 367 of a single object imposes global dependencies that constrain the textures at distant positions to be homogeneous (either B-left or B-center, but not B-right).

380 is not difficult to include it, and we do this in the  
 381 experiments presented later.) This gives

$$382 \Pr(I) \propto \sum_A \left[ \prod_l \prod_x \Pr(\mathbf{g}_l(x) | \mathbf{f}_{l+1}(x), A) \right] \Pr(A). \quad (4)$$

383 In principle, *any* distribution of images can be written in  
 384 this form since variables  $A$ , their joint distribution  $\Pr(A)$ ,  
 385 and the dependence of the image features on  $A$  can have  
 386 arbitrarily complex structure capturing any non-local behav-  
 387 ior.

388 Before we propose a specific structure for the variables  
 389  $A$  let us point out that the conditioning of  $\mathbf{g}_l(x)$  on  $\mathbf{f}_{l+1}(x)$   
 390 already captures some of the coarse-to-fine dependency of  
 391 image statistics. Many image structures, such as edges,  
 392 persist across scale, and so it is found that modeling this  
 393 dependency of features across scales is essential for the  
 394 synthesis of natural texture images (Portilla and Simoncelli,  
 395 2000). Note that we choose to condition  $\mathbf{G}_l$  on  $\mathbf{F}_{l+1}$   
 396 instead of  $\mathbf{G}_{l+1}$  as in Luetzgen and Willsky (1995). We  
 397 believe that this better captures local correlation since it is  
 398 consistent with empirically established natural image  
 399 statistics (Buccigrossi and Simoncelli, 1998), and with  
 400 equal image dimensions the conditioning becomes straight-  
 401 forward.

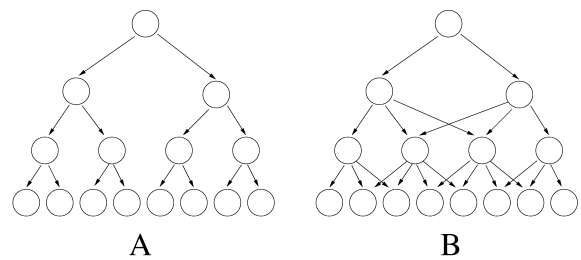
402 Eq. (4) can be seen as a mixture model with mixture  
 403 labels  $A$  conditioning the entire image. We remind the  
 404 reader that mixture models group samples with common  
 405 statistics by assigning them a common label (Duda et al.,  
 406 2001). In this case a sample corresponds to the entire  
 407 image. However, instead of entire images we intend to  
 408 group individual pixels in the pyramid. We consider,  
 409 therefore, the set of hidden variables as an unsupervised  
 410 segmentation. As such, we assign to each position and  
 411 layer in the pyramid a variable  $a_l(x)$  that conditions the  
 412 features only locally. This gives

$$413 \Pr(I) \propto \sum_A \left[ \prod_l \prod_x \Pr(\mathbf{g}_l(x) | \mathbf{f}_{l+1}(x), a_l(x)) \right] \Pr(A). \quad (5)$$

414 The structure of the joint distribution  $\Pr(A)$  of variables,  
 415  $A = \{a_l(x) | x \in \mathcal{P}_l, l = 0, \dots, L\}$ , captures the statistical  
 416 relation between the segmentation in different regions and  
 417 scales. In addition, due to the factorization over space the  
 418 dependency structure of  $\Pr(A)$  has to communicate non-  
 419 local information over different regions of the image and  
 420 across scale. A tree, as shown in Fig. 3(A), satisfies that  
 421 requirement and makes the necessary computations tractable.  
 422 With this choice the joint distribution is given by<sup>6</sup>

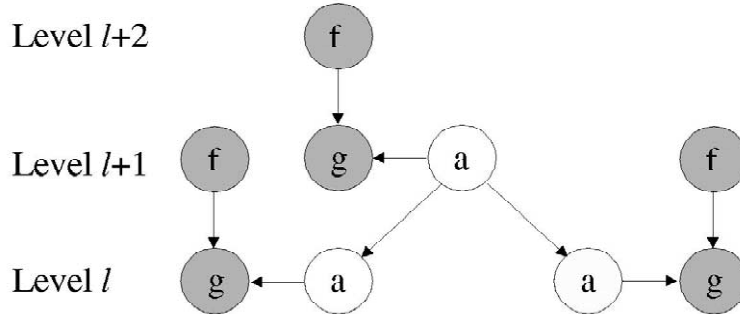
$$423 \Pr(A) = \prod_l \prod_x \Pr(a_l(x) | a_{l+1}(Px)), \quad (6)$$

424 where the probability  $\Pr(a_l(x) | a_{l+1}(Px))$  is that of finding  
 425  $a_l$  at  $x$  given  $a_{l+1}$  at the parent of  $x$ ,  $Px$ . We allow the  
 426 number of possible values for the labels  $a_l$  to be different



369  
 370 Fig. 3. Dependency structure for the label pyramid. (A) In a binary tree  
 371 probabilities can be propagated efficiently. The disadvantage is that some  
 372 neighboring nodes are very weakly linked, while others are very tightly  
 373 linked. (B) Dense graph where the smallest clique is the entire graph and  
 374 probability computations increase exponentially with the tree size.

375 <sup>6</sup>Variable  $a_{L+1}$  has not been defined and can be thought of as a label  
 376 with a single possible value. The conditional distribution  $\Pr(a_L | a_{L+1})$   
 377 then turns into a prior  $\Pr(a_L)$ . The reader should note that this footnote  
 378 applies to the remainder of the paper, most notably in the derivation of  
 379 the expectation in Section 3.2.



430 Fig. 4. Dependency structure of the HIP model corresponding to Eq. (7). To simplify the diagram, we show the dependency graph for a single parent node  
 431 conditioning two of its children. In practice, each parent has four children, i.e., a quadtree. Dark shaded nodes represent observable data. We also omit the  
 432 subscripts which indicate position. White nodes are hidden variables.

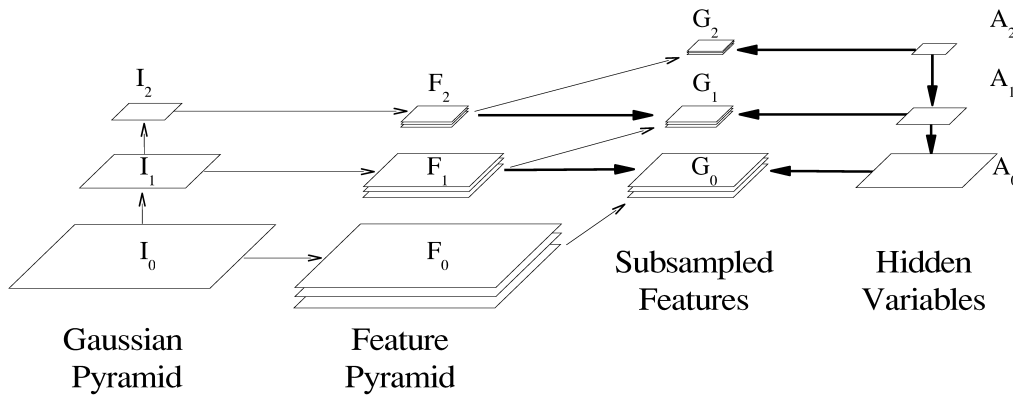
445 for each level  $l$ . Combining Eq. (6) with Eq. (5), and using  
 446 shorter notation<sup>7</sup> for the position  $x$ , we obtain

447 
$$\Pr(I) \propto \sum_A \prod_l \prod_x \Pr(\mathbf{g}_l | \mathbf{f}_{l+1}, a_l, x) \Pr(a_l | a_{l+1}, x). \quad (7)$$

448 The dependency graph of expression (7) is shown in Fig.  
 449 4. Note that this is not the only way to introduce hidden  
 450 variables to capture non-local dependence. The more  
 451 general model is still given by expression (4). However,  
 452 expression (7) represents a fairly general class of models  
 453 with several desired properties, in particular dependencies  
 454 proceeding from coarse-to-fine scales that are local in both  
 455 space and scale. The integration of the hidden variable  
 456 structure into the pyramid framework is depicted in Fig. 5.

### 3. Training the HIP model with an EM algorithm 457

We adjust the parameters of our model to match the 458  
 statistics of a given set of images by using Maximum 459  
 Likelihood (ML) parameter estimation. The structure of 460  
 the model in Eq. (7) and illustrated in Fig. 4 permits the 461  
 exact and efficient computation of all marginal prob- 462  
 abilities required for the expectation-maximization (EM) 463  
 algorithm (Dempster et al., 1977). The algorithm first 464  
 computes the expectations, over the hidden variables, of 465  
 the log-likelihood for a given set of parameters and 466  
 observations (E-step). Then, using these expectations, the 467  
 likelihood is maximized with respect to the parameters of 468  
 the model (M-step): 469



434  
 435 Fig. 5. The addition of hidden variable images for capturing long-range dependencies. Conditioning is shown with thick arrows while construction of  
 436 features is shown with thin arrows. In this example  $L = 2$ .

437 <sup>7</sup>In the following we will write  $\Pr(a_l(x) | a_{l+1}(Px))$  as simply  
 438  $\Pr(a_l | a_{l+1}, x)$ . For brevity we also write  $\Pr(\mathbf{g}_l | \mathbf{f}_{l+1}, a_l, x)$  for  
 439  $\Pr(\mathbf{g}_l(x) | \mathbf{f}_{l+1}(x), a_l(x))$ , the probability distribution for finding the feature  
 440 vector  $\mathbf{g}_l(x)$  at position  $x$  given that the feature vector  $\mathbf{f}_{l+1}(x)$  and hidden  
 441 variable  $a_l(x)$  were also found at  $x$ . Similar notation will be used for other  
 442 expressions. The argument  $x$  in  $\Pr(\cdot | \cdot, x)$  selects the random variables  
 443 associated with position  $x$  and should not be understood as a random  
 444 variable by itself.

E-step: 
$$Q(\theta | \theta^t) = \sum_A \Pr(A | I, \theta^t) \ln \Pr(I, A | \theta), \quad (8) \quad 470$$

M-step: 
$$\theta^{t+1} = \arg \max_{\theta} Q(\theta | \theta^t). \quad (9) \quad 471$$

Here we have summarized all parameters of the model in 472  
 $\theta$ , and  $\theta^t$  represents the values of the parameters in the 473  
 current iteration step  $t$ . 474

476 The main challenge for this model lies in computing the  
 477 expectations over the unknown labels. In this section, only  
 478 the resulting equations will be given. For the derivation of  
 479 the probability propagation in this hierarchical model  
 480 readers are referred to Appendix A.

### 481 3.1. Maximization

482 We start with the M-step by inserting Eq. (7) into Eq.  
 483 (8):

$$484 \quad Q(\theta | \theta^t) \\ 485 = \sum_A \Pr(A | I, \theta^t) \sum_l \sum_x \ln \Pr(\mathbf{g}_l, a_l | \mathbf{f}_{l+1}, a_{l+1}, x, \theta) + \text{const.}, \\ 486 \quad (10)$$

$$487 \\ 488 = \sum_l \sum_x \sum_{a_l, a_{l+1}} \{ \Pr(a_l, a_{l+1} | I, x, \theta^t) \ln \Pr(\mathbf{g}_l, a_l | \mathbf{f}_{l+1}, a_{l+1}, x) \} \\ 489 + \text{const.} \quad (11)$$

490 Here,  $\Pr(a_l, a_{l+1} | I, x, \theta^t)$  represents the marginal prob-  
 491 abilities of pairs of labels from neighboring layers at  
 492 position  $x$  for given image data and the current parameter  
 493 values. The additive constant is due to the proportionality  
 494 factors of Eq. (7). Assuming we know the probability  
 495  $\Pr(a_l, a_{l+1} | I, x, \theta^t)$  for all parent/child label pairs,  $a_l, a_{l+1}$ ,  
 496 we can search for the optimal parameters. At this point we  
 497 must commit to a parameterization of  $\Pr(a_l | a_{l+1}, x)$  and  
 498  $\Pr(\mathbf{g}_l | \mathbf{f}_{l+1}, a_{l+1}, x)$ . We use the same parameters for all  
 499 positions so that we obtain homogeneous behavior across  
 500 the image, a constraint which is often referred to as  
 501 “tying” in the HMM literature (Rabiner, 1989), and is also  
 502 used in the HMT model (Crouse et al., 1998). However,  
 503 we allow our model to have different parameters at  
 504 different pyramid levels—we tie across position but not  
 505 scale. We allow  $\Pr(a_l | a_{l+1}, x)$  to depend on the position of  
 506 the child relative to the parent, e.g. the probability is  
 507 different for the upper-right child than for the lower-left  
 508 child, etc. We also choose to parameterize  $\Pr(a_l | a_{l+1}, x)$  as

$$509 \quad \Pr(a_l | a_{l+1}) = \frac{\pi_{a_l, a_{l+1}}}{\sum_{a_l} \pi_{a_l, a_{l+1}}}. \quad (12)$$

510 The values of the parameters  $\pi_{a_l, a_{l+1}}$  are determined during  
 511 the updates in the M-step of the EM algorithm. Note that  
 512 with this definition  $\Pr(a_l | a_{l+1})$  is always properly normal-  
 513 ized. There is an arbitrary scale in the  $\pi_{a_l, a_{l+1}}$  for each  
 514  $a_{l+1}$ , but this is fixed by choosing a particular form for the  
 515 update equation. Note also that we omit  $x$  in this notation  
 516 as the parameterization is independent of the position  
 517 within a layer.

518 We assume a simple model for the distribution of  
 519 subsampled features conditioned on the features of the next  
 520 highest pyramid level. Our model represents a mixture  
 521 where the label  $a$  selects the mixture component. We

choose a Gaussian distribution where the parameters are  
 indexed by the labels and the dependency of the features is  
 parameterized as a linear relationship in the mean:

$$522 \quad \Pr(\mathbf{g} | \mathbf{f}, a) = \mathcal{N}(\mathbf{g}, M_a \mathbf{f} + \bar{\mathbf{g}}_a, \Lambda_a). \quad (13) \quad 525$$

526 If different features at a given spatial location (a pixel)  
 527 are independent, then diagonal  $M$  and  $\Lambda$  are sufficient. The  
 528 parameter set is now defined as

$$529 \quad \theta = \cup_l \{ \pi_{a_l, a_{l+1}}, M_{a_l}, \bar{\mathbf{g}}_{a_l}, \Lambda_{a_l} | a_l \in \{1, \dots, N_{a_l}\} \}.$$

530 With the choices (12) and (13) the M-step is easily  
 531 solved. The maximum of (11) with respect to  $\theta$  can be  
 532 found by setting the derivatives with respect to the  
 533 different parameters equal to zero and solving for the  
 534 corresponding parameter. For  $\pi_{a_l, a_{l+1}}^{t+1}$  we find

$$535 \quad \pi_{a_l, a_{l+1}}^{t+1} = \sum_x \Pr(a_l, a_{l+1} | I, x, \theta^t). \quad (14)$$

536 For the remaining update equations we define the  
 537 following weighted average:

$$538 \quad \langle X \rangle_{t, a_l} = \frac{\sum_x \Pr(a_l | I, x, \theta^t) X(x)}{\sum_x \Pr(a_l | I, x, \theta^t)}.$$

539 The weights  $\Pr(a_l | I, x, \theta^t)$  represent the marginal prob-  
 540 abilities of finding label value  $a_l(x)$  at position  $x$  given the  
 541 image data and the current parameter values.

The update equations are

$$542 \quad \bar{\mathbf{g}}_{a_l}^{t+1} = \langle \mathbf{g}_l \rangle_{t, a_l} - M_{a_l}^{t+1} \langle \mathbf{f}_{l+1} \rangle_{t, a_l}, \quad (16) \quad 543$$

$$544 \quad M_{a_l}^{t+1} = (\langle \mathbf{g}_l \mathbf{f}_{l+1}^T \rangle_{t, a_l} - \bar{\mathbf{g}}_{a_l}^{t+1} \langle \mathbf{f}_{l+1}^T \rangle_{t, a_l}) \times \langle \mathbf{f}_{l+1} \mathbf{f}_{l+1}^T \rangle_{t, a_l}^{-1}, \quad (17)$$

and

$$545 \quad \Lambda_{a_l}^{t+1} = \langle (\mathbf{g}_l - M_{a_l}^{t+1} \mathbf{f}_{l+1})(\mathbf{g}_l - M_{a_l}^{t+1} \mathbf{f}_{l+1})^T \rangle_{t, a_l} - \bar{\mathbf{g}}_{a_l}^{t+1} \bar{\mathbf{g}}_{a_l}^{t+1 T}. \quad (18) \quad 547$$

548 Since these expressions are mutually dependent, we must  
 549 insert Eq. (16) into Eq. (17) and solve for  $M_{a_l}^{t+1}$  to obtain

$$550 \quad M_{a_l}^{t+1} = (\langle \mathbf{g}_l \mathbf{f}_{l+1}^T \rangle_{t, a_l} - \langle \mathbf{g}_l \rangle_{t, a_l} \langle \mathbf{f}_{l+1}^T \rangle_{t, a_l}) (\langle \mathbf{f}_{l+1} \mathbf{f}_{l+1}^T \rangle_{t, a_l} \\ 551 - \langle \mathbf{f}_{l+1} \rangle_{t, a_l} \langle \mathbf{f}_{l+1}^T \rangle_{t, a_l})^{-1}.$$

To summarize, the update procedure is:

- 552 1. compute  $M_{a_l}^{t+1}$  according to Eq. (19),
- 553 2. compute  $\bar{\mathbf{g}}_{a_l}^{t+1}$  according to Eq. (16), then
- 554 3. compute  $\Lambda_{a_l}^{t+1}$  according to Eq. (18).

555 If we assume diagonal  $M$  and  $\Lambda$  we can ignore the  
 556 off-diagonal terms in these expressions. In fact, the com-  
 557 ponent densities  $\mathcal{N}(\mathbf{g}, M_a \mathbf{f} + \bar{\mathbf{g}}_a, \Lambda_a)$  factor into individual  
 558 densities for each component of  $\mathbf{g}$ . We can replace Eqs.  
 559 (19), (16), and (18) with their scalar versions and apply  
 560 them to each component of  $\mathbf{g}$  independently. 561

566 **3.2. Expectation**

567 In the E-step we compute the marginal probabilities of  
 568 pairs of labels from neighboring layers  $\Pr(a_l, a_{l+1} | I, x, \theta^t)$   
 569 for given image data. However, note that, in all occur-  
 570 rences of the re-estimation equations, i.e., Eqs. (12), (14)  
 571 and (15), we need that quantity only up to an overall  
 572 factor. We can choose that factor to be  $\Pr(I | \theta^t)$  and  
 573 therefore compute  $\Pr(a_l, a_{l+1}, I | x, \theta^t)$  using

574 
$$\Pr(a_l, a_{l+1} | I, x, \theta^t) \Pr(I | \theta^t) = \Pr(a_l, a_{l+1}, I | x, \theta^t)$$
  
 575 
$$= \sum_{A \setminus \{a_l(x), a_{l+1}(x)\}} \Pr(I, A | \theta^t). \quad (20)$$

576 The complexity of computing these sums relates to the  
 577 dependency structure of the variables  $A$ , which we have  
 578 already defined in Eq. (7) and Fig. 4.

579 From the viewpoint of computational complexity, it is  
 580 important to understand the rationale for this choice. From  
 581 the literature on graphical models (Jordan, 1998) we know  
 582 that the cost of evaluating these sums grows exponentially  
 583 with the clique sizes in the graph and linearly with the  
 584 number of cliques. If we choose the dependency such that  
 585 every label is conditioned on only one label from the  
 586 parent layer then the clique size is minimal (Fig. 3(A)). For  
 587 an image pyramid with subsampling-by-two that corre-  
 588 sponds to a quadtree structure. In a quadtree a location  $x_l$   
 589 has only one parent  $Px_l$  in layer  $l + 1$ , and four children  
 590  $Cx_l$  in layer  $l - 1$ . If we do not restrict the dependencies,  
 591 and maintain instead a more general belief network  
 592 structure between layers, with local connectivity (Fig.  
 593 3(B)), the entire label pyramid is one irreducible clique,  
 594 and the exact evaluation of the sums becomes prohibitive.

595 We now compute the probability of hidden labels given  
 596 the entire image pyramid. This computation will be  
 597 essentially the same as propagating the probabilities of  
 598 observations of the entire pyramid to a particular junction  
 599 of label pairs. Probabilities first propagate upwards, and  
 600 then downward to a particular label pair. During the  
 601 propagation we marginalize over the other labels. We  
 602 recursively define quantities  $u$  and  $d$ , representing the  
 603 upwards and downwards propagating probabilities:

604 
$$u_l(a_l, x) = \Pr(\mathbf{g}_l | \mathbf{f}_{l+1}, a_l, x) \prod_{x' \in Cx} \tilde{u}_{l-1}(a_l, x'), \quad (21)$$

$$\tilde{u}_l(a_{l+1}, x) = \sum_{a_l} \Pr(a_l | a_{l+1}) u_l(a_l, x), \quad (22) \quad 605$$

$$d_l(a_l, x) = \sum_{a_{l+1}} \Pr(a_l | a_{l+1}) \tilde{d}_l(a_{l+1}, x), \quad (23) \quad 606$$

$$\tilde{d}_l(a_{l+1}, x) = \frac{u_{l+1}(a_{l+1}, Px)}{\tilde{u}_l(a_{l+1}, x)} d_{l+1}(a_{l+1}, Px). \quad (24) \quad 607$$

The upward recursion (Eqs. (21) and (22)) is initialized at  
 608  $l = 0$  with  $u_0(a_0, x) = \Pr(\mathbf{g} | \mathbf{f}_1, a_0, x)$  and ends at  $l = L$ . At  
 609 layer  $L$ , Eq. (22) reduces to  $\tilde{u}_L(a_{L+1}, x) = \tilde{u}_L(x)$ . Since we  
 610 do not model any further dependencies beyond layer  $L$ , the  
 611 pixels at layer  $L$  are assumed independent. The product of  
 612 all  $\tilde{u}_L(x)$  is the total image probability:  
 613

$$\Pr(I | \theta^t) = \prod_{x \in \mathcal{P}_L} \tilde{u}_L(x) = u_{L+1}. \quad (25) \quad 614$$

The downward recursion (Eqs. (23) and (24)) starts with  
 615 Eq. (24) at  $l = L$  with  $d_{L+1}(a_{L+1}, x) = d_{L+1}(x) = 1$ , and  
 616 ends at  $l = 0$  with Eq. (23).  
 617

618 With these quantities we can compute Eq. (20) as

$$\Pr(a_l, a_{l+1}, I | x, \theta^t) = u_l(a_l, x) \tilde{d}_l(a_{l+1}, x) \Pr(a_l | a_{l+1}), \quad (26) \quad 619$$

$$\Pr(a_l, I | x, \theta^t) = u_l(a_l, x) d_l(a_l, x), \quad (27) \quad 620$$

621 where the computations (21)–(27) in the E-step at iteration  
 622  $t$  are performed with fixed parameters  $\theta^t$ .

623 **3.3. Emission probabilities**

624 The model described thus far uses the same labels  $A$  for  
 625 modeling the distributions of the observables  $\mathbf{G}$  as well as  
 626 for propagating non-local information through the different  
 627 scales. For the latter purpose it might be necessary to have  
 628 many different possible label values that can encode for  
 629 more complex information. In the levels of the pyramid the  
 630 means and variances that are assigned to each label value  
 631 may have very few pixels for training and therefore may  
 632 be poorly estimated. It is thus reasonable to separate the  
 633 functionality of the label  $A$ , for example as indicated in  
 634 Fig. 6. In this case, labels  $A$  still code for the non-local  
 635 information while labels  $B$  now are used for modeling the  
 636 distribution of the features. Up to the conditioning on  $\mathbf{F}$

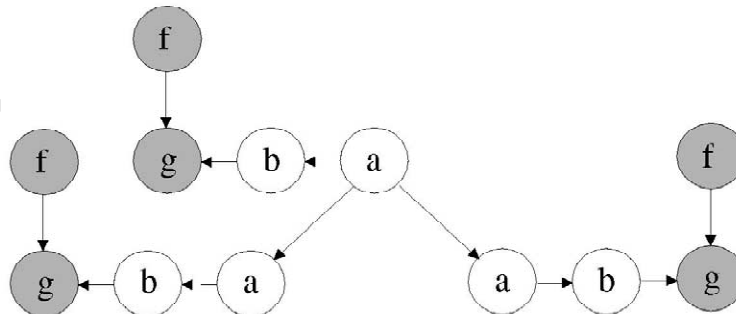


Fig. 6. Dependency structure of the HIP model with emission probabilities corresponding to Eq. (28).



638 this model now very closely resembles an HMM tree, with  
639 mixture densities as emission probabilities.

640 The expressions for the joint probability distributions as  
641 well as the corresponding probability propagation can be  
642 obtained by setting

$$643 \Pr(\mathbf{g}_l | \mathbf{f}_{l+1}, a_l, x) = \sum_{b_l} \Pr(b_l | a_l) \Pr(\mathbf{g}_l | \mathbf{f}_{l+1}, b_l, x) \quad (28)$$

644 in Eqs. (7) and (21).

645 Once again we parameterize  $\Pr(b_l | a_l)$  as

$$646 \Pr(b_l | a_l) = \frac{\pi_{b_l, a_l}}{\sum_{b_l} \pi_{b_l, a_l}}. \quad (29)$$

647 The re-estimation equation in the M-step is then

$$648 \pi_{b_l, a_l}^{t+1} = \sum_x \Pr(b_l, a_l | I, x, \theta^t), \quad (30)$$

649 where we can use the joint

$$650 \Pr(b_l, a_l, I | x, \theta^t) = \frac{\Pr(\mathbf{g}_l | \mathbf{f}_{l+1}, b_l, x)}{\Pr(\mathbf{g}_l | \mathbf{f}_{l+1}, a_l, x)} \Pr(b_l | a_l) \Pr(a_l, I | x, \theta^t). \quad (31)$$

### 652 3.4. Scale, mixture and hierarchy labels

653 An alternative to adding separate labels for mixture  
654 components as emission probabilities is to partition the  
655 labels in the simple model, in other words to view a label  
656  $a$ , as in Fig. 4, as being composed of a label  $m \in$   
657  $\{1, \dots, N_m\}$  that specifies the mixture component and a  
658 hierarchy label  $c \in \{1, \dots, N_c\}$  that is intended to capture  
659 non-local information. We can relate these labels to each  
660 other in different ways, for example  $a = (c - 1)N_m + m$   
661 while requiring  $N_a = N_m N_c$ . The mixture component at a  
662 location in the pyramid is given by  $m$ , whereas  $c$  influences  
663 image structure at finer scales through the model's condi-  
664 tional probability distribution  $\Pr(a_l | a_{l+1})$ . We can now  
665 choose a small value for  $N_m$  at low-resolution levels, and a  
666 larger value for  $N_c$ . Conversely, the only appropriate value  
667 for  $N_c$  at the finest-resolution level is one, since all  
668 information from other levels can be carried to  $m_l$  by  $a_{l+1}$ .

669 We can recover emission probabilities from this model  
670 by imposing a simplification, namely that  $\Pr(a_l | a_{l+1}) =$   
671  $\Pr(m_l | c_l) \Pr(c_l | c_{l+1})$ . In this case,  $N_c$  at the finest-res-  
672 olution level should be greater than one.

673 We can go further and add more structure to the labels.  
674 In the models we use for mass detection, we further  
675 partition the mixture labels into a label  $m$  and a *scale* label  
676  $z$ , so that

$$677 \Pr(\mathbf{g}_l | \mathbf{f}_{l+1}, m_l, z_l) = \mathcal{N}(\mathbf{g}_l, M_{m_l} \mathbf{f}_{l+1} + \sigma_{z_l} \bar{\mathbf{g}}_{m_l}, \sigma_{z_l}^2 \Lambda_{m_l}). \quad (32)$$

678 The mixture components for a given value of  $m_l$  model the  
679 same type of image structure, but with means and vari-  
680 ances of different magnitudes as determined by the scale  
681 factor  $\sigma_{z_l}$ . Such explicit representation of scale has been

682 reported to be important in modeling natural image  
683 structure (Wainwright and Simoncelli, 1999; Wainwright et  
684 al., 2001; Romberg et al., 2001). As with these earlier  
685 models, we make the scale factors depend on their parents  
686 at lower-resolution levels, since the magnitude of wavelet  
687 coefficients tends to persist across pyramid level. To  
688 reduce the number of model parameters we chose to  
689 constrain the label probabilities so that

$$690 \Pr(a_l | a_{l+1}) = \Pr(m_l | c_{l+1}) \Pr(z_l | z_{l+1}, c_{l+1}) \Pr(c_l | c_{l+1}). \quad (33)$$

692 Note that there is ambiguity in this representation, since we  
693 can multiply  $\bar{\mathbf{g}}_{m_l}$  by a factor  $\lambda$  and  $\Lambda_{m_l}$  by  $\lambda^2$  for all  $m_l$ ,  
694 and the mixture components will not change if we also  
695 multiply  $\sigma_{z_l}$  by  $\lambda^{-1}$  for all  $z_l$ . To remove the ambiguity we  
696 apply the constraint  $\prod_{z_l} \sigma_{z_l} = 1$ .

697 The M-step of the EM algorithm must be modified in  
698 this model, since we cannot solve for all of the parameters  
699 of  $\Pr(\mathbf{g}_l | \mathbf{f}_{l+1}, m_l, z_l)$  analytically. We choose to use a  
700 generalized EM algorithm, in which the M-step is iterative.  
701 Since the E-step is more computationally intensive than the  
702 M-step, the increase in training time is relatively small. In  
703 the iterative M-step, we fix  $\sigma_{z_l}$  for all  $z_l$  and re-estimate the  
704 other mixture parameters to maximize the expected likeli-  
705 hood. We then fix these other parameters (that depend only  
706 on  $m_l$ ) and re-estimate  $\sigma_{z_l}$  for all  $z_l$ . Within the M-step we  
707 alternate repeating these two substeps several times. In  
708 practice, two to four iterations usually is adequate.

709 Finally, we allow for rotations into the label structure,  
710 although the use of wavelets restricts us to rotations by  
711 multiples of  $90^\circ$ . We do this simply by requiring that for  
712 every value of  $m_l$ , there are three other values whose mean  
713 and covariance ( $\bar{\mathbf{g}}_{m_l}$  and  $\Lambda_{m_l}$ ) are related to those of the  
714 original by the three rotations of  $90^\circ$ ,  $180^\circ$ , and  $270^\circ$ . The  
715 EM algorithm is easily modified to handle this. In the  
716 E-step we can compute expectations for each individual  
717 mixture component. For a given set of four components  
718 related by rotations, we first perform the appropriate  
719 inverse rotation on each of three of the expectation values,  
720 so they agree with the fourth component, and sum them. In  
721 the M-step we then update the parameters for this com-  
722 ponent from these sums, and copy the results to the other  
723 three components after rotating them, as appropriate.

### 724 3.5. Preprocessing and training methods

725 We divide the data set into training and test sets of  
726 approximately equal size and, for the mass detection, we  
727 construct a jackknife (i.e. 10 different random splits) so as  
728 to demonstrate the robustness of the results. We use a set of  
729 approximately orthogonal wavelets to decompose the  
730 intensity images into feature images (see Appendix B for  
731 details). Before applying the wavelet decomposition we  
732 wrap images at edges in order to obtain perfect reconstruc-

tion for the compression and synthesis. We crop images so that they are square with objects approximately centered.

We train the HIP model using the EM algorithm described in Section 3. The number of labels was chosen through a splitting procedure, using the minimum description length (MDL) cost criterion to compute the optimal model. The MDL cost is given as  $-\log \Pr(I | H) + d/2 \log(N)$  (Rissanen, 1996; Deco and Obradovic, 1996), where  $\Pr(I | H)$  is the probability of the training data under the model  $H$ ,  $d$  is the number of parameters in  $H$ , and  $N$  is the number of images in the training set (note that the form of the HIP model makes it well suited for MDL model selection). For the splitting procedure, we begin with only one hidden label value at each level. We then duplicate each label along with its parameters, i.e.,  $\Pr(a_l | a_{l+1})$ ,  $\bar{g}_{a_l}$ ,  $M_{a_l}$ , and  $A_{a_l}$ , randomly perturbing the duplicate parameters. We re-train the new, larger model and compare its MDL cost with the previous model. We repeat this, successively duplicating labels, retraining, and evaluating the MDL cost, until the MDL cost increases. The model with the lowest MDL cost is then used in the applications presented below. Fig. 7 shows how the area under the ROC curve  $A_z$  for the test data tracks the MDL cost. Note that best performance tends to be for lowest costs, though the tracking is not monotonic.

Such an MDL-based training procedure is feasible, however it is computationally expensive. On a Sun Ultrasparc-2 workstation the entire splitting procedure required roughly two weeks of computer time. This is partly due to the large number of parameters being adjusted, 12,995 for the optimal model for the masses (see Appendix C for discussion on the number of parameters in the model). In spite of this, over-fitting does not appear to

be a problem, as evidenced from our jackknife results presented below. We believe this is so because every location in each level of the wavelet decomposition provides examples for the parameters used to model the coefficients at that level. For example, as part of fitting the image distribution the model must fit the marginal distributions of the wavelet coefficients at each level. For this purpose there is one example per location in each image, so that there are many more effective examples than the number of images. Also, once the models are trained, there is minimal computational cost/overhead in applying them for detection, synthesis and compression.

#### 4. Experimental results

In this section we report results for applying a HIP model, with complete scale, hierarchy and mixture labels, to mammographic image analysis, in particular detection of mammographic masses. As an experimental paradigm, we choose to demonstrate the utility of the approach for a dataset representing the output of the University of Chicago's CAD system for mass detection. This is a state-of-the-art mammographic screening system which includes a set of signal enhancement, pre-processing, rule-based and statistical-based classification schemes for detecting masses in digitized mammograms (Doi et al., 1993; Nishikawa et al., 1996; Giger et al., 2000). We choose this

<sup>8</sup>For example, the Digital Database for Screening Mammography (DDSM), the Mammographic Image Analysis Society (MIAS) database, and the Lawrence Livermore National Laboratories (LLNL)/University of California at San Francisco (UCSF) database.

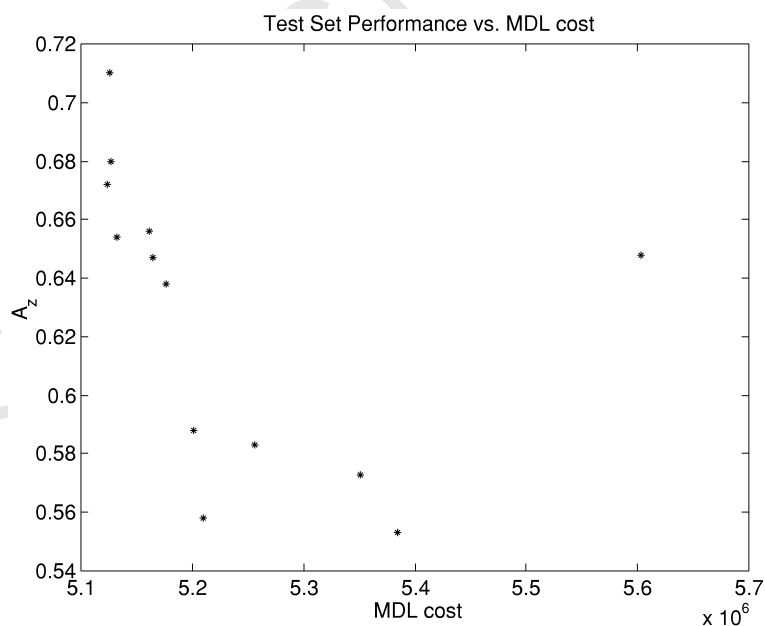


Fig. 7. Performance (area under the ROC curve, or  $A_z$ ) versus MDL cost for the HIP model pair used to generate the experimental results.

814 paradigm over an alternative, such as performance on a  
 815 public database of digitized mammograms,<sup>8</sup> since we can  
 816 better estimate the clinical impact of the model in terms of  
 817 reducing difficult false positives as well as demonstrating  
 818 performance relative to a well-characterized clinical sys-  
 819 tem. As an additional demonstration of the utility of HIP,  
 820 we compare results to that of an HMT model using a  
 821 single set of hidden labels to model two component  
 822 mixtures over a wavelet tree. Details of this model can be  
 823 found in Crouse et al. (1998) and Romberg et al. (2001).  
 824 The comparison with the HMT enables us to characterize  
 825 performance relative to another hierarchical probabilistic  
 826 model for images, specifically the utility of the additional  
 827 hidden label structure.

828 In the following we first describe the dataset used in the  
 829 experiments and then present our results for classification,  
 830 synthesis and compression.

#### 831 4.1. Mammographic dataset

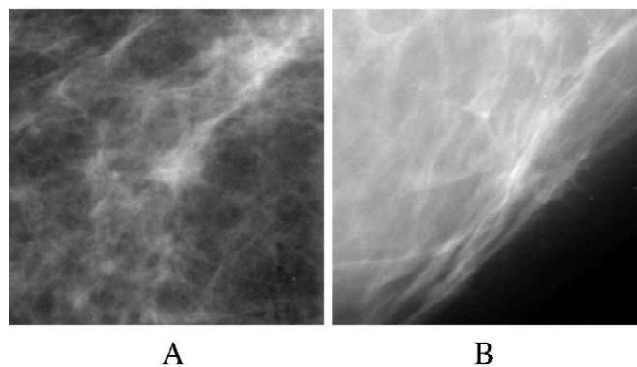
832 The dataset used in these experiments consists of  
 833 regions of interest (ROIs) selected from digitized mam-  
 834 mograms by the mammographic mass CAD system devel-  
 835 oped by the Rossmann Laboratories of the University of  
 836 Chicago (UofC) (Nishikawa et al., 1995, 1996; Giger et  
 837 al., 2000). The CAD system, consisting of a series of  
 838 classification/detection modules, places an indicator (e.g.,  
 839 “arrow”) next to potential masses on a digital image of  
 840 the mammogram. The location of the indicator is deter-  
 841 mined by dividing the mammogram into ROIs and then  
 842 eliminating false-positive ROIs using pattern recognition  
 843 techniques. The output of the CAD system therefore can be  
 844 seen as a set of ROIs, of which all are assumed to be  
 845 positive for masses. Since this is a screening system, both  
 846 malignant and benign masses are considered “true posi-  
 847 tives”. ROIs output by the CAD system which do not  
 848 contain masses (i.e. non-masses) are UofC false positives.

849 For the experiments in this paper, 169 ROIs were  
 850 available, of which 72 contained masses (true positives)  
 851 and 97 were false positives of the UofC CAD system. The  
 852 detected objects (apparent masses) are not necessarily  
 853 centered in the ROI, since they may lie close to the edge of  
 854 the mammogram. The original ROIs are 512-by-512  
 855 pixels. Examples of mass and non-mass ROIs are shown in  
 856 Fig. 8.

#### 857 4.2. Mass detection

858 We first consider using HIP as a post-processor (i.e.  
 859 adjunct) to the UofC CAD system (Nishikawa et al.,  
 860 1996). The goal was to determine if the HIP model could  
 861 be used to reduce false positives without reducing sen-  
 862 sitivity. In addition, the performance of the HIP model was  
 863 compared to an HMT. A 10-way jackknife was used to  
 864 compute the results.

865 Two HIP models were trained for each of the jackknife



801  
 802 Fig. 8. Examples of data used in experiments. (A) Mammographic mass  
 803 (true positive). (B) False positive generated by the UofC CAD system.

866 sets. Each jackknife set consisted of 36 randomly chosen  
 867 ROIs that contained masses, and 48 randomly chosen ROIs  
 868 without masses. One model was trained for the mass ROIs  
 869 and another model for the non-mass ROIs. Similarly, two  
 870 HMTs were trained, using the same split of the data.

871 The likelihood ratio under the two models was used as  
 872 the test criterion, i.e., a threshold on this ratio is used to  
 873 decide which ROIs will be detected as masses. The true  
 874 and false-positive fractions as a function of the threshold  
 875 were measured on the split of the jackknife that contained  
 876 the test set. This set also consisted of 36 mass and 49  
 877 non-mass ROIs.

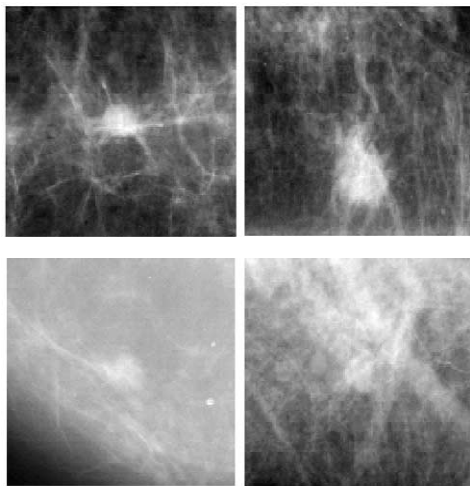
878 Table 1 summarizes the results for the jackknife experi-  
 879 ments. On average, the receiver operating characteristic  
 880 (ROC) curve (Metz, 1988) for the HIP model applied to  
 881 the test images has an area under the curve ( $A_z$ ) equal to  
 882 0.78 and a 16% reduction in false positives generated by  
 883 the UofC CAD system, without loss in sensitivity. By  
 884 comparison the HMT model has a mean  $A_z$  equal to 0.55,  
 885 with only a 3% reduction in false positives. Given the  
 886 difficulty of this dataset (i.e. it represents the most difficult  
 887 false positives that could not be eliminated by the UofC  
 888 system) a 16% reduction in the false-positive rate is  
 889 significant. Nonetheless, the HIP model is not capable of  
 890 learning subtle differences for distinguishing between  
 891 masses and non-masses. Fig. 9 shows examples of ROIs  
 892 correctly and incorrectly classified by HIP. From these  
 893 examples we see that the model trained to detect masses  
 894 performs well for ROIs containing localized, and some-  
 895 what isolated, homogeneous “mass-like” structure. For  
 896 non-masses (UofC false positives) the HIP model correctly  
 897 characterizes ROIs devoid of mass-like structure, and in

804 Table 1  
 805 Jackknife results for mass detection  
 806

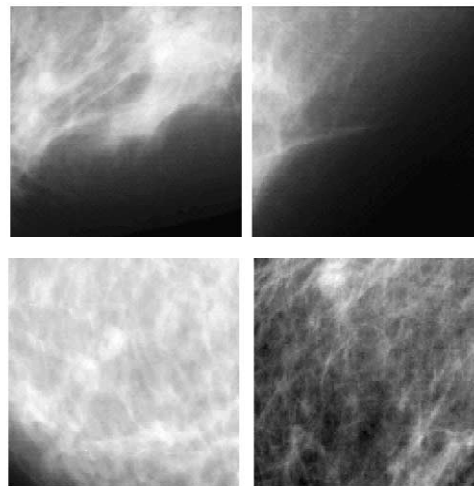
	HIP	HMT
Mean (std) $A_z$	0.78 (0.04)	0.55 (0.05)
Mean (std) FPF@100% TPF	0.84 (0.15)	0.97 (0.05)
Mean (std) FPF@95% TPF	0.73 (0.17)	0.93 (0.05)

811 FPF, false-positive fraction; TPF, true-positive fraction.  
 812

masses



non-masses (UofC fps)



900

901 Fig. 9. Example of ROIs which the HIP model correctly (top row) and incorrectly (bottom row) characterizes. Note that the difference between the two  
 902 classes of ROI (mass versus non-mass) is much more apparent in the top row than in the bottom row, consistent with model performance.

903 fact learns that many of the non-mass false positives are in  
 904 fact at the breast border. For ROIs incorrectly character-  
 905 ized by HIP, we see a striking similarity in the ROI  
 906 structure for masses and non-masses, which is consistent  
 907 with the fact that the HIP model would have difficulty  
 908 putting these ROIs into one of the two classes. More  
 909 insight into the structure captured by the mass and non-  
 910 mass HIP models, and how it relates to detection per-  
 911 formance, can be seen through mammographic synthesis.

#### 912 4.3. Mammographic synthesis

913 Since the HIP model is a generative model, we can  
 914 sample the model and synthesize new images. In the  
 915 context of ROI classification, synthesized images can  
 916 provide qualitative insight into what features the model is  
 917 extracting and representing for both positive and negative  
 918 ROIs. Using the same ROI database used for classification,  
 919 we constructed HIP models for positives (masses) and  
 920 negatives (no masses). The trained HIP models were  
 921 sampled to synthesize new ROI images. The sampling  
 922 procedure begins at the coarsest resolution, where the  
 923 hidden labels are randomly sampled from the distribution  
 924  $\Pr(A_L)$ . The feature images  $\mathbf{G}_L$  are then sampled from  
 925  $\Pr(\mathbf{G}_L | A_L)$ . The  $\mathbf{G}_L$  are used to construct  $I_{L-1}$ , from  
 926 which the  $\mathbf{F}_L$  are constructed. We then sample  $A_{L-1}$  from  
 927  $\Pr(A_{L-1} | A_L)$ , and then  $\mathbf{G}_{L-1}$  from  $\Pr(\mathbf{G}_{L-1} | \mathbf{F}_L, A_{L-1})$ .  
 928 This is repeated until the finest resolution is reached and  $I_0$   
 929 is constructed.

930 Fig. 10 shows examples of these images. Inspection of  
 931 the synthesized positive ROIs shows more focal structure,  
 932 with more well-defined borders and higher spatial fre-  
 933 quency content than the negative ROIs. Comparison to the  
 934 HMT synthesized images, constructed with a similar  
 935 sampling procedure, shows the HMT images for positive

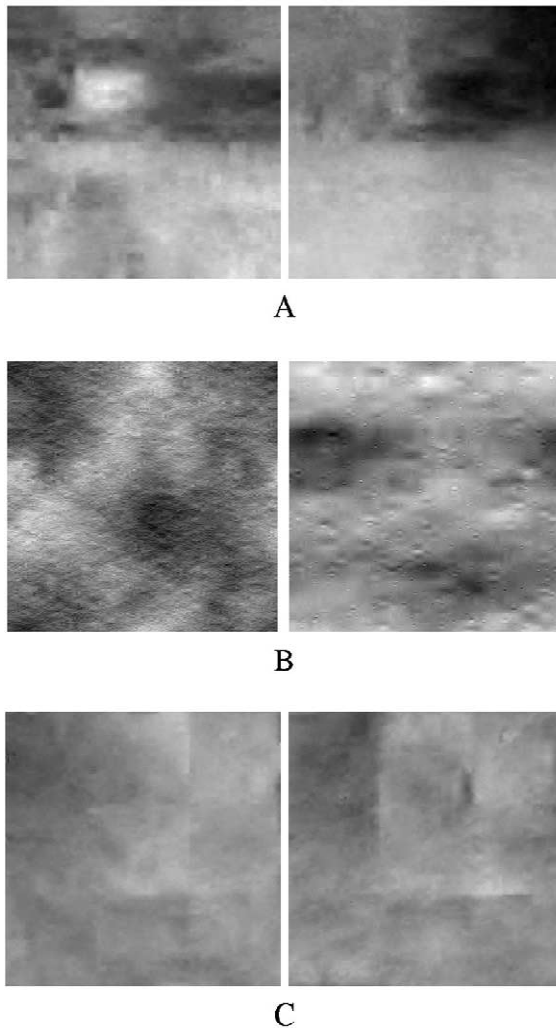
and negative ROIs. Though positive and negative ROIs are  
 936 different, the positive ROI does not capture the focal  
 937 structure of the mass, as is the case for the HIP generated  
 938 images. This is likely due to the flexibility of the hidden  
 939 variable architecture of the HIP, with scale, mixture and  
 940 hierarchy labels able to capture more structure in the  
 941 image. As a test, we sampled a HIP model constructed  
 942 using only a single hidden label structure (similar to that of  
 943 an HMT). Fig. 10(C) shows that the focal structure has  
 944 now disappeared in the positive ROI.  
 945

946 It is equally important to consider the mammographic  
 947 structure that is not well represented in the synthesized  
 948 images. A comparison of Fig. 8(A) and Fig. 10(A) indi-  
 949 cates that the model is not accurately representing the  
 950 extended linear structure of the breast parenchyma. One  
 951 possible reason is that the tree structure of the model is not  
 952 ideal for capturing colinear dependencies across space,  
 953 since there is no direct conditioning between neighboring  
 954 nodes. Such dependencies can be captured only indirectly  
 955 via propagation up and down the tree.

#### 956 4.4. Mammographic image compression

957 A stream of random variables can be optimally com-  
 958 pressed if we know their distribution. A HIP model of a  
 959 source of images should therefore allow us to compress  
 960 examples of those images with high efficiency. Here we  
 961 demonstrate compression with HIP and HMT models using  
 962 a simple technique.

963 Given an image and a HIP model, we compress the  
 964 image as follows. First, we compute the most likely value  
 965 of each hidden label,  $a_l^*(x) = \arg \max_{a_l} \Pr(a_l, I | x, \theta^l)$ ,  
 966 using Eq. (27). These most likely values are then encoded  
 967 with arithmetic coders, which require a probability dis-  
 968 tribution for the symbols they are to encode. For this we



971

972 Fig. 10. Mammographic ROI images synthesized from positive and  
 973 negative HIP and HMT models. (A) Synthesized ROIs from HIP model  
 974 with scale, hierarchy and mixture labels. Positive ROIs (left) tend to have  
 975 more focal structure, with more defined borders and higher spatial  
 976 frequency content. Negative ROIs (right) tend to be more amorphous with  
 977 lower spatial frequency content. (B) Synthesized ROIs from HMT model.  
 978 Though ROIs of positive and negative models appear different, the  
 979 positive ROI does not appear to capture the focal structure of masses. (C)  
 980 HIP model with single label architecture. As with the HMT, this  
 981 architecture does not capture the focal structure of the masses.

982 use the HIP model distributions  $\Pr(a_i^*(x) | a_{i+1}^*(x))$ . Given  
 983 the label value  $a_i^*(x)$ , we then encode the feature vector  
 984  $\mathbf{g}_i(x)$  using  $\Pr(\mathbf{g}_i | \mathbf{f}_{i+1}, a_i^*, x)$ . The latter is used by de-  
 985 composing  $\mathbf{g}_i(x)$  into its components along the eigenvec-  
 986 tors of the covariance matrix of  $\Pr(\mathbf{g}_i | \mathbf{f}_{i+1}, a_i^*)$ ,  $\Lambda_{a_i^*}$ .  
 987 These components are independent under  $\Pr(\mathbf{g}_i | \mathbf{f}_{i+1}, a_i^*)$ ,  
 988 so they can be encoded independently. Each component is  
 989 encoded with a specified precision by dividing the real line  
 990 into intervals of width equal to twice the precision. Using  
 991 an arithmetic coder we then encode an index for the  
 992 interval containing the component. The probability of each  
 993 bin is provided by the integral of the univariate Gaussian  
 994 distribution of the component implied by  $\Pr(\mathbf{g}_i | \mathbf{f}_{i+1}, a_i^*)$ ,

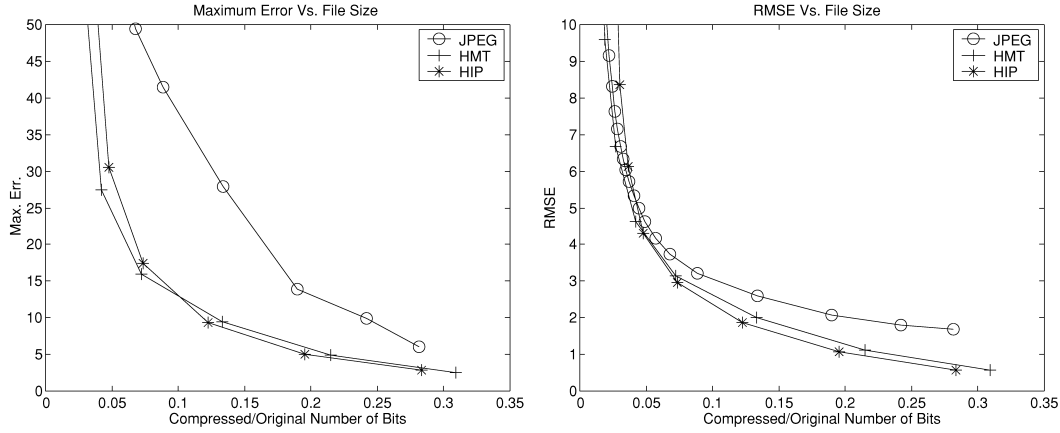
i.e., with variance given by the corresponding eigenvalue  
 of  $\Lambda_{a_i^*}$ . This procedure is computationally expensive, and  
 is not necessarily optimal even if the HIP model exactly  
 matches the image distribution, but it serves to demon-  
 strate the capability. We used the analogous procedure for  
 compression with the HMT models, except the DC re-  
 sidual needs to be stored separately, without compression.

We compress the images at several values of the  
 precision. The range of precisions was chosen to roughly  
 match the errors given by JPEG with a range of quality  
 factors.

To compare with JPEG, we first convert the images to  
 unsigned byte pixels. We divide the pixel values by four  
 before compressing with JPEG, since the maximum value  
 of the pixels in the mass ROIs is a little less than 1024, and  
 multiply by four after decompressing. The results, aver-  
 aged across all images, are shown in Fig. 11. The HIP  
 model performs better than HMT for higher precision or  
 lower loss, suggesting that the hidden labels capture useful  
 information, allowing better compression. The better per-  
 formance of HMT at higher compression ratios comes  
 from our method for encoding the hidden labels. In our  
 scheme these use the same number of bits no matter what  
 the precision is, putting a lower limit on the compressed  
 image size. The HMT model has fewer and simpler hidden  
 values to be encoded. It still has a lower limit, but this is  
 much smaller than for the HIP model. In fact, the HMT  
 model performs better than JPEG at these high compres-  
 sion ratios. A more sophisticated compression algorithm  
 with the HIP model would group labels, since some  
 mixture components can become indistinguishable when  
 coding at low precision. This grouping would effectively  
 adjust the complexity of the HIP model with coding  
 precision.

## 5. Discussion and conclusion

We have developed a class of multi-scale probabilistic  
 network models for images which we call hierarchical  
 image probability or HIP models. To justify these, we  
 show that image distributions can be exactly represented as  
 products over pyramid levels of distributions of sub-sam-  
 pled feature images conditioned on coarser-scale image  
 information. We argue that hidden variables are needed to  
 capture long-range dependencies while allowing us to  
 further factor the distributions over position. In our current  
 model the hidden variables include scale, hierarchy and  
 mixture labels which enable a more flexible modeling of  
 natural images, compared to the structure of an HMT. This  
 was demonstrated by comparison of the two approaches  
 for mammographic mass detection, synthesis and compres-  
 sion, with the HIP model giving superior results. However,  
 the current structure of HIP is not well suited for capturing  
 dependencies between oriented spatial structure. Future  
 work will investigate methods for more direct modeling of



1050

1051 Fig. 11. Pixel error vs. size of compressed files for JPEG, HIP, and HMT. (Left) Maximum error ( $L_\infty$  norm). (Right) RMS error. These curves represent  
1052 averages across all of the mammographic ROIs.

1053 spatial orientation dependencies, which are obvious in the  
1054 structure of mammograms and, in general, natural images.

1055 Because HIP models are probability distributions over  
1056 images, they can be used for a wide range of image  
1057 processing tasks, e.g. classification, compression, noise  
1058 suppression, up-sampling, error correction, etc. In fact, any  
1059 image analysis problem can be approached in a principled  
1060 way using such distributions. Here we have presented  
1061 results for mammographic image analysis. However, there  
1062 are obviously other modalities and medical application  
1063 areas where HIP models would be useful. One in particular  
1064 is multi-modal fusion, where the problem is to bring a set  
1065 of images, acquired using different imaging modalities,  
1066 into alignment. One method that has demonstrated partic-  
1067 ularly good performance uses mutual information as an  
1068 objective criterion (Wells et al., 1996). The computation of  
1069 mutual information requires an estimate of entropies,  
1070 which in turn requires an estimate of the underlying  
1071 densities of the images. The HIP model potentially pro-  
1072 vides a framework for learning those densities.

1073 Some of the results we have obtained with the HIP  
1074 model are comparable to those given by other approaches  
1075 rather than being superior to them (e.g., for detection the  
1076 HPNN (Sajda et al., 2002) gives similar if not better  
1077 results). However, we obtain our results for several differ-  
1078 ent problems using a single model, rather than training  
1079 very different models for each problem. This flexibility  
1080 and the principled approach provided by HIP models to  
1081 image analysis are quite useful. We believe that, with  
1082 further development, models of image probability dis-  
1083 tributions will give superior performance in a variety of  
1084 medical image processing tasks.

1085 **Acknowledgements**

1086 We thank Robert Nishikawa and Maryellen Giger for  
1087 useful discussions and providing the data and Adam

Gerson for assistance in the simulations. We also thank the  
three anonymous reviewers for their helpful comments  
which greatly improved the manuscript. This work was  
funded by the U.S. Army Medical Research and Material  
Command (DAMD17-98-1-8061). P.S. was also supported  
by the DoD Multidisciplinary University Research Initia-  
tive (MURI) program administered by the Office of Naval  
Research under grant N00014-01-1-0625 as well as a grant  
from the National Imagery and Mapping Agency,  
NMA201-02-C0012.

1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097

1098 **Appendix A. Belief propagation in HIP**

1099 Here we show how to obtain the upwards and down-  
1100 wards propagation rules (21)–(24). All the computations  
1101 can be executed locally. Consider the subgraph presented  
1102 in Fig. A.1. In this subgraph, every node  $X$  can take on a  
1103 discrete number of values, with  $\sum_X$  indicating a sum over  
1104 those values. Assigned to every node  $X$  is also an evidence  
1105 node  $g_X$ , with a fixed value for given image data.  $g_X \dots$   
1106 refers to  $g_X$  and all the evidence in the rest of the graph  
1107 that can be reached through node  $X$ . Using this notation the  
1108 entire evidence provided by the image  $I$  is the collection  
1109  $\{g_A \dots, g_B \dots, g_C \dots\}$ . The probability required in the EM  
1110 algorithm is

1111 
$$\Pr(B, A, I) = \Pr(B, A, g_A \dots, g_B \dots, g_C \dots), \quad (A.1)$$

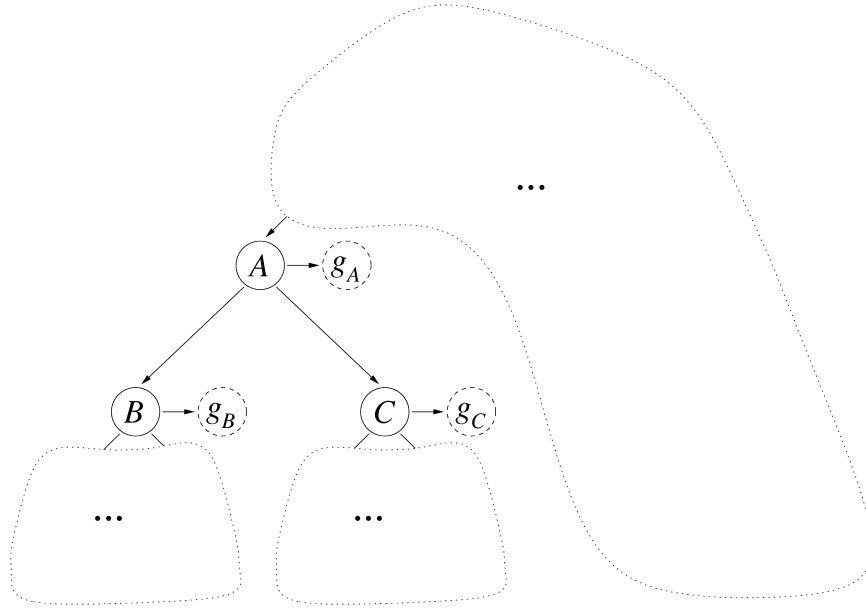
1112 
$$= \Pr(A, g_A \dots, g_C \dots) \Pr(B, g_B \dots | A), \quad (A.2)$$

1113 
$$= \Pr(A, g_A \dots, g_C \dots) \Pr(B | A) \Pr(g_B \dots | B), \quad (A.3)$$

1114 
$$= d_B(A) \Pr(B | A) u(B), \quad (A.4)$$

1115 where in Eq. (A.4) the quantities  $d_B(A)$  and  $u(B)$  represent  
1116

1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116



1119

1120 Fig. A.1. Subgraph of the label pyramid. Conditioned on  $A$  the variables that are connected to  $A$  become independent, such as labels  $B$ ,  $C$ , and the  
 1121 evidence node  $g_A$ . These variables are also conditionally independent to the joint variables that can be reached going upwards to the rest of the tree  
 1122 structure.

1123 the probabilistic influences on node  $B$  that are down-  
 1124 stream and up-stream from the node. In (A.2) we use the  
 1125 fact that conditioned on  $A$  the evidence coming through the  
 1126 children of  $A$  is independent from the rest of the tree  
 1127 beyond  $A$ . Since the children of  $A$  have no other parents,  
 1128 all the probabilistic influence beyond that parent edge can  
 1129 be communicated only through  $A$ . Similarly in (A.3) we  
 1130 use the fact that the evidence  $g_B$  is independent from the  
 1131 children of  $B$  if conditioned on  $B$ . Finally, in (A.4) we use  
 1132 the definitions for computing these probabilities recursive-  
 1133 ly in an upwards and downwards probability propagation  
 1134 as follows:

$$1135 \quad u(A) \equiv \Pr(g_A, g_B \dots, g_C \dots \mid A), \quad (A.5)$$

$$1136 \quad = \Pr(g_A \mid A) \Pr(g_B \dots \mid A) \Pr(g_C \dots \mid A), \quad (A.6)$$

$$1137 \quad = \Pr(g_A \mid A) u_B(A) u_C(A)$$

$$1138 \quad = \Pr(g_A \mid A) \prod_{x \in \mathcal{C}_A} u_x(A), \quad (A.7)$$

$$1139 \quad u_B(A) \equiv \Pr(g_B \dots \mid A), \quad (A.8)$$

$$1140 \quad = \sum_B \Pr(B \mid A) \Pr(g_B \dots \mid B), \quad (A.9)$$

$$1141 \quad = \sum_B \Pr(B \mid A) u(B). \quad (A.10)$$

1142 We use in (A.6) and (A.9) conditional independence  
 1143 when conditioning on  $A$  and  $B$ , respectively. In (A.10) we  
 1144 use definition (A.5) for node  $B$  and in (A.7) we use  
 1145 definition (A.8) for the children of  $A$ . The downward  
 1146 propagating probability is defined and computed as

$$d_B(A) = \Pr(A, g_A \dots, g_C \dots), \quad (A.11) \quad 1147$$

$$= \Pr(g_C \dots \mid A) \Pr(A, g_A \dots), \quad (A.12) \quad 1148$$

$$= \frac{u(A)}{u_B(A)} d(A), \quad (A.13) \quad 1149$$

$$d(B) \equiv \Pr(B, g_A \dots, g_C \dots), \quad (A.14) \quad 1150$$

$$= \sum_A \Pr(B \mid A) \Pr(A, g_A \dots, g_C \dots), \quad (A.15) \quad 1151$$

$$= \sum_A \Pr(B \mid A) d_B(A). \quad (A.16) \quad 1152$$

Again, we use the conditional independences when  
 1153 conditioning on  $A$  in (A.12), (A.13), and (A.15). One can  
 1154 verify (A.13) by inserting the corresponding definitions  
 1155 and canceling the term  $\Pr(g_A \mid A)$  to recover (A.12).  
 1156

These upwards and downwards propagation equations  
 1157 are the basis for Eqs. (21)–(24).  
 1158

## Appendix B. Wavelets 1159

For the HIP model presented in this paper we use  
 1160 approximately orthogonal wavelets with subsampling by  
 1161 two. The two-dimensional filters are separable, being  
 1162 products of one-dimensional wavelets. For the one-dimen-  
 1163 sional filters we solve for appropriate tap weights (i.e. filter  
 1164 coefficients) subject to the following constraints:  
 1165

1. One filter is even-symmetric and low-pass (taps sum to  
 1167 one, zero response at the Nyquist frequency). 1167
2. The second filter is odd-symmetric (high-pass). 1168

1170 Table B.1

1171 Tap weights for 12-tap orthogonal wavelet filters with subsampling by two

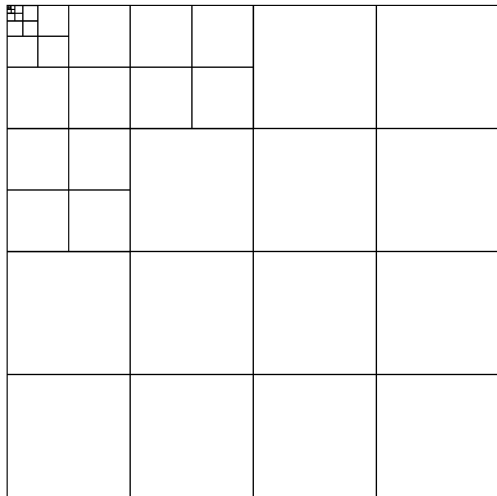
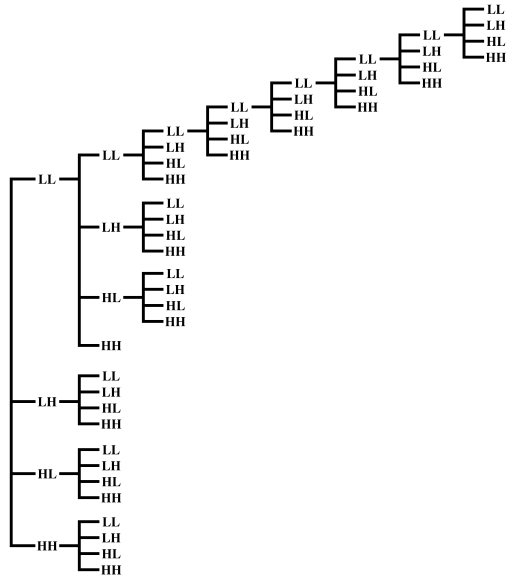
1173 Tap index	0	1	2	3	4	5
1174 Tap weights	0.492631	0.060246	-0.060468	0.004588	0.002779	0.000223

1175

- 1212 3. The first derivatives of the responses of the filters are  
 1213 zero at zero frequency and the Nyquist frequency.  
 1214 4. The second derivatives of the responses of the filters are  
 1215 zero at zero frequency and the Nyquist frequency.  
 1216 We adjusted the tap weights to be as close to orthogonal  
 1217 as possible, by minimizing the mean-squared error of

analysis and subsequent reconstruction of a noise image. 1218  
 No number of taps gave perfect reconstruction, and we 1219  
 decided that the error with 12 taps was acceptable. Note 1220  
 that the odd-symmetric high-pass filter has the same tap 1221  
 weights as the low-pass filter, except for an alternating 1222  
 sign. Numerical values for the tap weights of the filters are 1223  
 given in Table B.1. 1224

For the HIP model we build a wavelet packet tree using 1225  
 the entropy minimization techniques developed by Saito 1226  
 (1994). We use half the data to compute a wavelet packet 1227  
 with minimal entropy, which is analogous to maximizing 1228  
 the sparsity of the wavelet coefficients. The wavelet packet 1229  
 that is constructed is shown in Fig. A.2. Note that from this 1230  
 representation one can see the dimensionality of  $\mathbf{g}_l$  at each 1231  
 scale. 1232



### Appendix C. Parameters 1233

Each mixture component (label  $m$ ) has parameters  $\bar{\mathbf{g}}$ ,  $\Lambda$ , 1234  
 and  $M$ . If  $\mathbf{g}$  has dimension  $N_g$  and  $\mathbf{f}$  has dimension  $N_f$ , then 1235  
 $\bar{\mathbf{g}}$  is  $N_g$  parameters,  $\Lambda$  is  $N_g(N_g + 1)/2$  parameters, and  $M$  1236  
 is  $N_g N_f$  parameters. At level 2, for example,  $N_g = 13$  and 1237  
 $N_f = 12$ , so  $\bar{\mathbf{g}}$  is 13 parameters,  $\Lambda$  is 91 parameters, and  $M$  1238  
 is 156 parameters, for a total of 260 parameters per value 1239  
 of the label  $m$ . These values for all levels in the MDL 1240

Table C.1 1183

Parameter counts for mixture components in HIP mass model 1184  
 1185

Level	$N_g$	$N_f$	$\Lambda$	$M$	Per comp.	No. comps.	Total
2	13	12	91	156	260	16	4160
3	11	4	66	44	121	16	1936
4	3	4	6	12	21	16	336
5	3	4	6	12	21	16	336
6	3	4	6	12	21	4	84
7	3	4	6	12	21	4	84
8	4	0	10	0	14	1	14
Total							6950

1186  
1187  
1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197

Table C.2 1198

Parameter counts for scale components in HIP mass model 1199  
 1200

Level	$N_c$	Total
2	8	7
3	8	7
4	8	7
5	8	7
6	8	7
7	4	3
8	2	1
Total		39

1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211

1177  
 1178 Fig. A.2. Wavelet packet constructed using entropy minimization. (A)  
 1179 Tree structure of wavelet packet with minimum entropy found using  
 1180 mammographic mass data. (B) Corresponding wavelet packet decomposi-  
 1181 tion using conventional representation for images. All blocks of equal size  
 1182 make up the features at a given level,  $\mathbf{g}_l$ .



1242 Table C.3  
 1243 Parameter counts for conditional probability distributions in HIP mass model  
 1244

1245 1246	Level	$N_c$	$N_m$	$N_z$	$\Pr(c_i   c_{i+1})$	$\Pr(m_i   c_{i+1})$	$\Pr(z_i   z_{i+1}, c_{i+1})$	Total
1247	2	1	16	8	0	2016	1792	3808
1248	3	8	16	8	224	2016	1792	4032
1249	4	8	16	8	224	2016	1792	4032
1250	5	8	16	8	224	2016	1792	4032
1251	6	8	4	8	224	480	896	1600
1252	7	8	4	4	112	240	96	448
1253	8	4	1	2	48	48	0	96
1254	9	4	0	0	4	0	0	4
1255 1256	Total							18,052

1257 optimal positive mass model are shown in Table C.1, along  
 1258 with totals. Table C.2 gives the number of parameters for  
 1259 the scale component label values  $z$ . Note that, because of  
 1260 the sum-to-one constraint on  $z$ , there is one fewer parame-  
 1261 ter than  $N_z$ .

1262 The number of parameters in conditional probability  
 1263 distribution for labels is less than the product of the  
 1264 numbers of labels appearing in the distribution. This is  
 1265 partly due to the normalization, which reduces the count by  
 1266 one. In addition, for purposes of MDL model selection, we  
 1267 can argue that only non-zero probabilities should be  
 1268 counted. For coding purposes we could code a distribution  
 1269  $\Pr(a | b)$  for a given  $b$  as the number of values of  $a$  for  
 1270 which the distribution is non-zero, followed by the values  
 1271 for which it is non-zero, followed by the corresponding  
 1272 non-zero probabilities (except the last, which is given by  
 1273 the normalization). Asymptotically, as the precision with  
 1274 which we encode the probabilities increases, the cost of  
 1275 coding the integers becomes negligible because they are  
 1276 fixed, whereas the code length for the probabilities in-  
 1277 creases. In Table C.3 we list the maximum possible  
 1278 number of parameters for each level, that is, assuming no  
 1279 zero probabilities. Note that the number of mixture com-  
 1280 ponents is increased by a factor of four compared to Table  
 1281 C.1, due to the rotational symmetry we have imposed. (A  
 1282 set of four components related by rotations are specified by  
 1283 the same parameters, but we do not impose constraints on  
 1284 the label probabilities.) Also, since each non-leaf node in  
 1285 the tree has four children the number of parameters for one  
 1286 of the distributions is four times the number of labels at  
 1287 one level times the number of labels at the parent level.  
 1288 The exception is level 9, where there are four values of the  
 1289 label  $c_9$  that condition the one child at level eight. Of the  
 1290 18,052 parameters, 12,275 are zero.

1291 We arrive at the total number of parameters in the model  
 1292 as the number of mixture components (6950) added to the  
 1293 number of scale labels (39) added to non-zero parameters  
 1294 for the conditional probabilities (18,052 – 12,046 = 6006)  
 1295 to get 12,995. Since we jackknife our data this represents a  
 1296 rough estimate on the number of parameters. There are a  
 1297 similar number of parameters for the non-mass models.

## References

Bird, R., 1990. Professional quality assurance for mammography screen-  
 ing programs. *Radiology* 177, 8–10.

Buccigrossi, R.W., Simoncelli, E.P., 1998. Image compression via joint  
 statistical characterization in the wavelet domain. *Tech. Rep. 414*, U.  
 Penn. GRASP Laboratory, available at ftp://ftp.cis.upenn.edu/pub/  
 eero/buccigrossi97.ps.gz.

Burhenne, L., Wood, S., D’Orsi, C., Feig, S., Kopans, D., O’Shaughnessy,  
 K., Sickles, E., Tabar, L., Vyborny, C., Castellino, R., 2000. Potential  
 contribution of computer-aided detection to the sensitivity of screening  
 mammography. *Radiology* 215, 554–562.

Burt, P.J., Adelson, E.H., 1983. The Laplacian pyramid as a compact  
 image code. *IEEE Trans. Commun. COM-31* (4), 532–540.

Chan, H., Sahiner, B., Lam, K.L., Petrick, N., Helvie, M.A., Goodsitt, M.,  
 Adler, D., 1998. Computerized analysis of mammographic microcalci-  
 fications in morphological and feature spaces. *Med. Phys.* 25, 2007–  
 2019.

Chellappa, R., Chatterjee, S., 1985. Classification of textures using  
 Gaussian Markov random fields. *IEEE Trans. ASSP* 33, 959–963.

Cheng, H., Bouman, C.A., 2001. Multiscale Bayesian segmentation using  
 a trainable context model. *IEEE Trans. Image Process.* 10 (4), 511–  
 525.

Coi, H., Baraniuk, R.G., 2001. Multiscale image segmentation using  
 wavelet-domain hidden Markov models. *IEEE Trans. Image Process.*  
 10 (9), 1309–1321.

Cootes, T., Hill, A., Taylor, C., Haslam, J., 1994. The use of active shape  
 models for locating structure in medical images. *Image Vis. Comput.*  
 12 (6), 355–366.

Cootes, T., Taylor, C., 2001. Statistical models of appearance for medical  
 image analysis and computer vision. In: Sonka, M., Hanson, K. (Eds.).  
*Medical Imaging 2001*, Vol. 4322. SPIE Press, pp. 236–248.

Cover, T.M., Thomas, J.A., 1991. *Elements of Information Theory*.  
 Wiley, New York.

Crouse, M.S., Nowak, R.D., Baraniuk, R.G., 1998. Wavelet-based statisti-  
 cal signal processing using hidden Markov models. *IEEE Trans. Signal*  
*Process.* 46 (4), 886–902.

Dayan, P., Abbott, L., 2002. *Theoretical Neuroscience: Computational*  
*and Mathematical Modeling of Neural Systems*. MIT Press, Cam-  
 bridge, MA.

De Bonet, J.S., Viola, P., 1998. Texture recognition using a non-paramet-  
 ric multi-scale statistical model. In: *Conference on Computer Vision*  
*and Pattern Recognition*. IEEE, pp. 641–647.

De Bonet, J.S., Viola, P., Fisher, J.W., 1998. Flexible histograms: a  
 multiresolution target discrimination model. In: *Zelnic, E.G. (Ed.)*.  
*Proceedings of SPIE*, Vol. 3371, pp. 519–530.

Deco, G., Obradovic, D., 1996. *An Information-Theoretic Approach to*  
*Neural Computing*. Springer, New York.

Dempster, N.M., Laird, A., Rubin, D.B., 1977. Maximum likelihood from

1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345

- 1347 incomplete data via the EM algorithm. *J. R. Statist. Soc. B* 39, 1348 185–197.
- 1349 Doi, K., Giger, M., Nishikawa, R., Hoffmann, K., MacMahon, H.,  
1350 Schmidt, R., Chua, K., 1993. Digital radiography: a useful clinical tool  
1351 for computer-aided diagnosis by quantitative analysis of radiographic  
1352 images. *Acta Radiol.* 34, 426–439.
- 1353 Duda, R., Hart, P., Stork, D., 2001. *Pattern Classification*, 2nd Edition.  
1354 Wiley, New York.
- 1355 Floyd, C., Lo, J., Yun, A., Sullivan, D., Kornguth, P., 1994. Prediction of  
1356 breast cancer malignancy using an artificial neural network. *Cancer* 74,  
1357 2944–2948.
- 1358 Freeman, W.T., Jones, T.R., Pasztor, E.C., 2002. Example-based super-  
1359 resolution. *IEEE Comput. Graph. Applic.* 22 (2), 56–65.
- 1360 Geman, S., Geman, D., 1984. Stochastic relaxation, Gibbs distributions,  
1361 and the Bayesian restoration of images. *IEEE Trans. PAMI* 6 (6),  
1362 194–207.
- 1363 Giger, M., Huo, Z., Kupinski, M., Vyborny, C., 2000. Computer-aided  
1364 diagnosis in mammography. In: Sonka, M., Fitzpatrick, J. (Eds.),  
1365 *Medical Image Processing and Analysis. Handbook of Medical*  
1366 *Imaging*, Vol. 2. SPIE Press, pp. 917–986.
- 1367 Grenander, U., 1983. *Tutorials in Pattern Synthesis*. Brown University,  
1368 Providence, RI.
- 1369 Grenander, U., Chow, Y., Keenan, D., 1991. *Hands: A Pattern Theoretic*  
1370 *Study of Biological Shapes*. Springer, New York.
- 1371 Huo, Z., Giger, M., Vyborny, C., Wolverton, D., Schmidt, R., Doi, K.,  
1372 1998. Automated computerized classification of malignant and benign  
1373 mass lesions on digital mammograms. *Acad. Radiol.* 5, 155–168.
- 1374 Jiang, Y., Nishikawa, R., Wolverton, D., Metz, C., Giger, M.L., Schmidt,  
1375 R., Doi, K., 1996. Automated feature analysis and classification of  
1376 malignant and benign microcalcifications. *Radiology* 198, 671–678.
- 1377 Jordan, M.I. (Ed.), 1998. *Learning in Graphical Models*. NATO Science  
1378 Series D: Behavioral and Brain Sciences, Vol. 89. Kluwer Academic.
- 1379 Kopans, D., 1989. *Breast Imaging*. Lippincott, Philadelphia, PA.
- 1380 Lo, J., Kim, J., Baker, J., Floyd, C., 1996. Computer-aided diagnosis of  
1381 mammography using an artificial neural network: predicting the  
1382 invasiveness of breast cancers from image features. In: Giger, M.L.  
1383 (Ed.). *Medical Imaging 1996: Image Processing*, Vol. 2710. SPIE  
1384 Press, pp. 725–732.
- 1385 Luetzgen, M.R., Willsky, A.S., 1995. Likelihood calculation for a class of  
1386 multiscale stochastic models, with application to texture discrimina-  
1387 tion. *IEEE Trans. Image Proc.* 4 (2), 194–207.
- 1388 Metz, C., 1988. Current problems in ROC analysis. In: *Proceedings of the*  
1389 *Chest Imaging Conference, Madison, WI*, pp. 315–333.
- 1390 Metz, C., Shen, J., 1992. Gains in accuracy from replicated readings of  
1391 diagnostic images: prediction and assessment in terms of ROC  
1392 analysis. *Med. Decis. Making* 12, 60–75.
- Nishikawa, R., Schmidt, R., Osnis, R., Giger, M., Doi, K., Wolverton, D., 1393  
1996. Two-year evaluation of a prototype clinical mammographic 1394  
workstation for computer-aided diagnosis. *Radiology* 201, 256. 1395
- Nishikawa, R.M., Haldemann, R.C., Papaioannou, J., Giger, M.L., Lu, P., 1396  
Schmidt, R.A., Wolverton, D.E., Bick, U., Doi, K., 1995. Initial 1397  
experience with a prototype clinical intelligent mammography work- 1398  
station for computer-aided diagnosis. In: Loew, M.H., Hanson, K.M. 1399  
(Eds.). *Medical Imaging 1995*, Vol. 2434. SPIE, Bellingham, WA, pp. 1400  
65–71. 1401
- Portilla, J., Simoncelli, E., 2000. A parametric texture model based on 1402  
joint statistics of complex wavelet coefficients. *Int. J. Comput. Vis.* 40 1403  
(1), 49–71. 1404
- Rabiner, L., 1989. A tutorial on hidden Markov models and selected 1405  
applications in speech recognition. *Proc. IEEE* 77 (2), 257–285. 1406
- Rissanen, J., 1996. Information theory and neural nets. In: Smolensky, 1407  
Mozer, Rumelhart (Eds.), *Mathematical Perspectives on Neural Net-* 1408  
*works*. pp. 567–602. 1409
- Romberg, J.K., Coi, H., Baraniuk, R.G., 2001. Bayesian tree-structured 1410  
image modeling using wavelet domain hidden Markov models. *IEEE* 1411  
*Trans. Image Process.* 10 (7), 1056–1068. 1412
- Saito, N., 1994. Local feature extraction and its applications using a 1413  
library of bases. Tech. Rep., Ph.D. Thesis (Advisor: Prof. R.R. Coif- 1414  
man), Department of Mathematics, Yale University. 1415
- Sajda, P., Spence, C., Pearson, J., 2002. Learning contextual relationship 1416  
in mammograms using a hierarchical pyramid neural network. *IEEE* 1417  
*Trans. Med. Imaging* 21 (3), 239–250. 1418
- Thurfjell, E., Lernevall, K., Taube, A., 1994. Benefit of independent 1419  
double reading in a population-based mammography screening pro- 1420  
gram. *Radiology* 191, 241–244. 1421
- Wainwright, M., Simoncelli, E., Willsky, A., 2001. Random cascades on 1422  
wavelet trees and their use in analyzing and modeling natural images. 1423  
*Appl. Comput. Harmonic Anal.* 11, 89–123. 1424
- Wainwright, M.J., Simoncelli, E.P., 1999. Scale mixtures of Gaussians 1425  
and the statistics of natural images. In: Solla, S.A., Leen, T., Müller, 1426  
K.-R. (Eds.). *Advances in Neural Information Processing Systems*, Vol. 1427  
12. MIT Press, Cambridge, MA, pp. 855–861. 1428
- Wells, W.M., Viola, P., Atsumi, H., Nakajima, S., Kikinis, R., 1996. 1429  
Multi-modal volume registration by maximization of mutual infor- 1430  
mation. *Med. Image Anal.* 1 (1), 35–51. 1431
- Zhang, W., Doi, K., Giger, M.L., Wu, Y., Nishikawa, R.M., Schmidt, R., 1432  
1994. Computerized detection of clustered microcalcifications in 1433  
digital mammograms using a shift-invariant artificial neural network. 1434  
*Med. Phys.* 21 (4), 517–524. 1435
- Zhu, S.C., Wu, Y.N., Mumford, D., 1997. Minimax entropy principle and 1436  
its application to texture modeling. *Neural Comput.* 9 (8), 1627–1660. 1437