

# Graphical Models, Inference and Learning

Brendan J. Frey  
University of Toronto

Tutorial paper available at:

[www.psi.toronto.edu/~frey/stuff/tutorial.ps.gz](http://www.psi.toronto.edu/~frey/stuff/tutorial.ps.gz)

# Acknowledgements

- Nebojsa Jojic, Microsoft Research
- Resources:
  - Neal and Hinton 93, A new view of the EM algorithm
  - Jordan 98 (ed), Inference and Learning in Graphical Models
  - Wiberg, Loeliger, Koetter 95, The sum-product algorithm for error-correcting decoding
  - Neal 93, Probabilistic inference using Markov chain Monte Carlo techniques
  - ... (see tutorial paper)

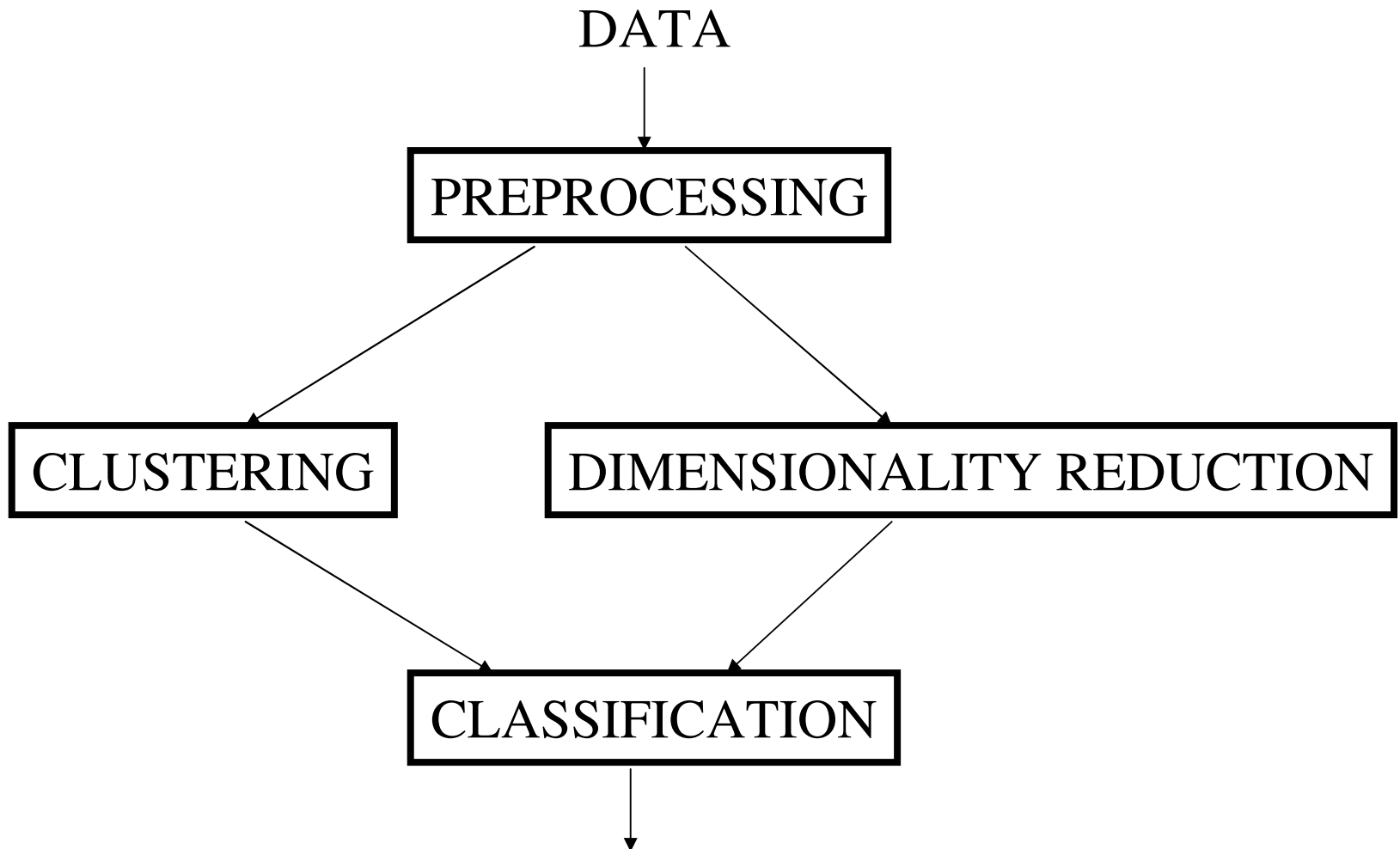
# Motivation

- Computer vision
- Communications systems
- Molecular biology
- Microphone array processing
- ...

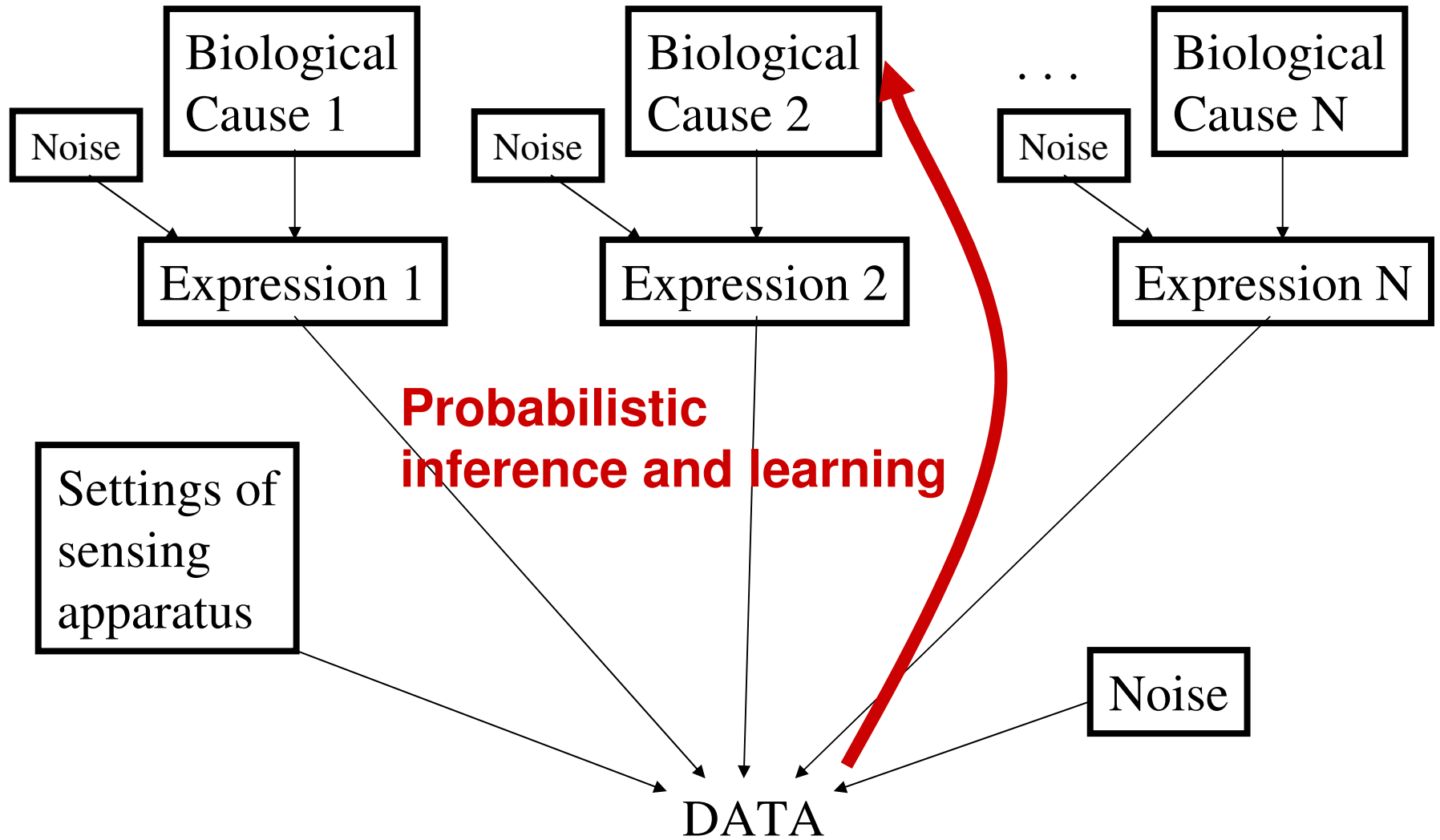
# Discriminative models and Generative models

- Discriminative:  $P(\text{class}|\text{input})$ 
  - Supervised learning
- Generative:  $P(\text{class}, \text{input})$ 
  - Supervised or unsupervised learning (class unobserved)
- Hidden variables:
  - $P(\text{class}|\text{input}) = \sum_{\text{hidden}} P(\text{class}, \text{hidden}|\text{input})$
  - $P(\text{class}, \text{input}) = \sum_{\text{hidden}} P(\text{class}, \text{hidden}, \text{input})$

# Sequential (block-diagram) data analysis



# The generative approach (on a genomics problem)



# Modularity and Graphical Models

- $P(\text{class}, \text{hidden}, \text{input})$  is complex for real-world problems
  - How do we specify constraints on variables?
  - How do we modify sub-components of the model?
  - How do we cope with computational intractability?
  - How do we cope with learnability?
- Graphical models: Modular descriptions of complex probability models

# Case study: Occlusion Model

B. J. FREY & N. JOJIC



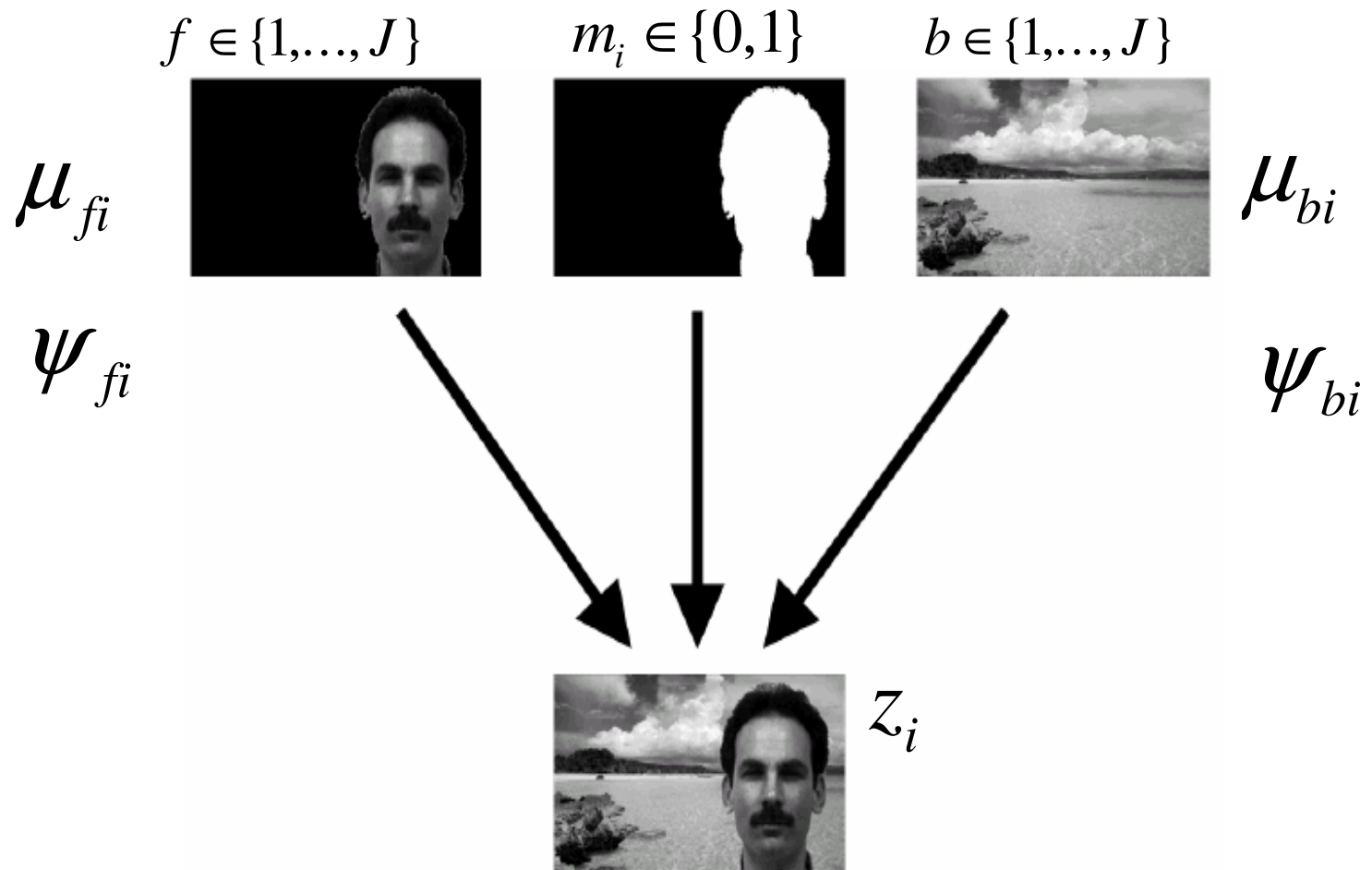
# A sample from the dataset

(5 different faces against 7 different backgrounds)

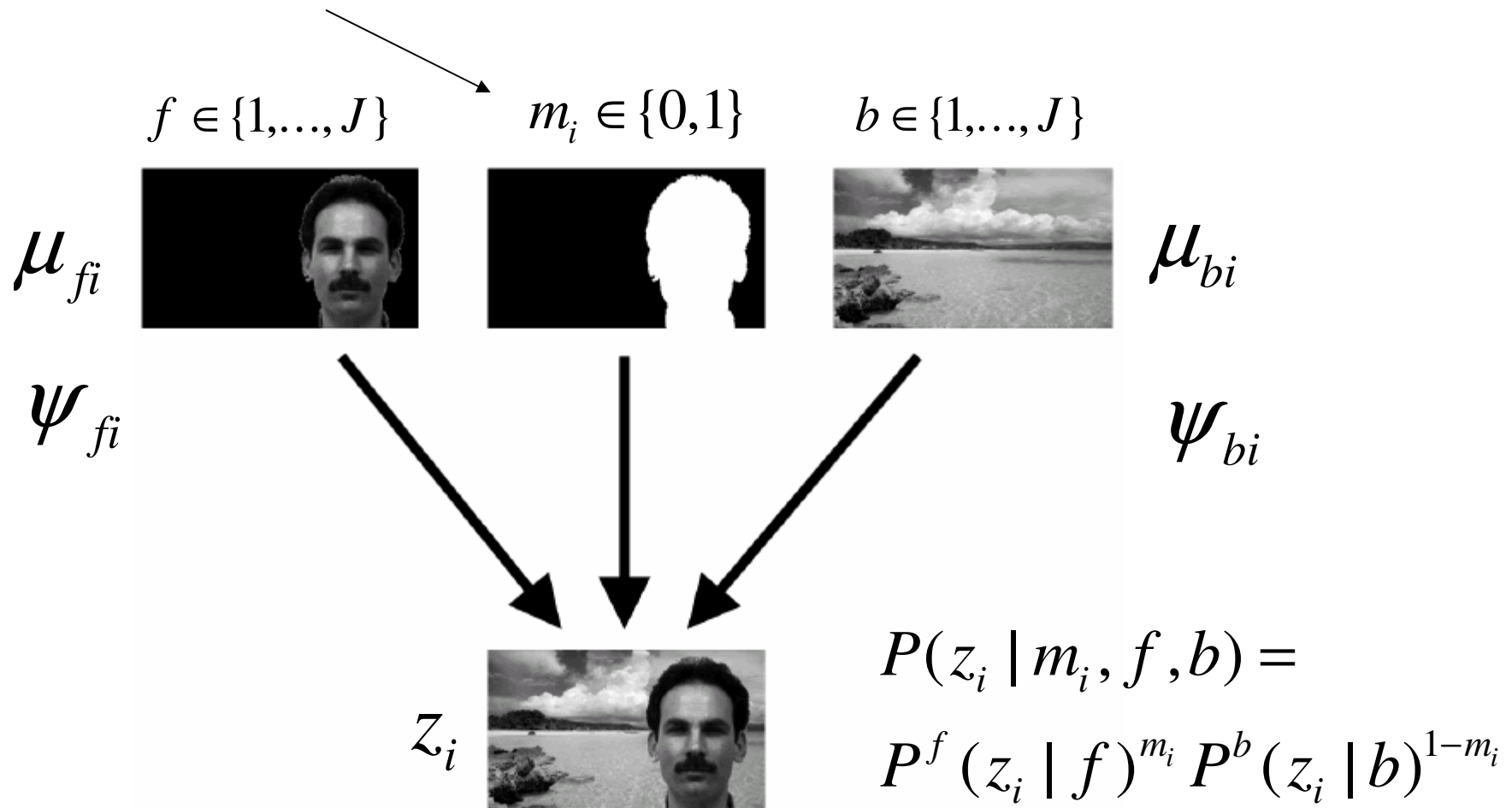


B. J. FREY & N. JOJIC

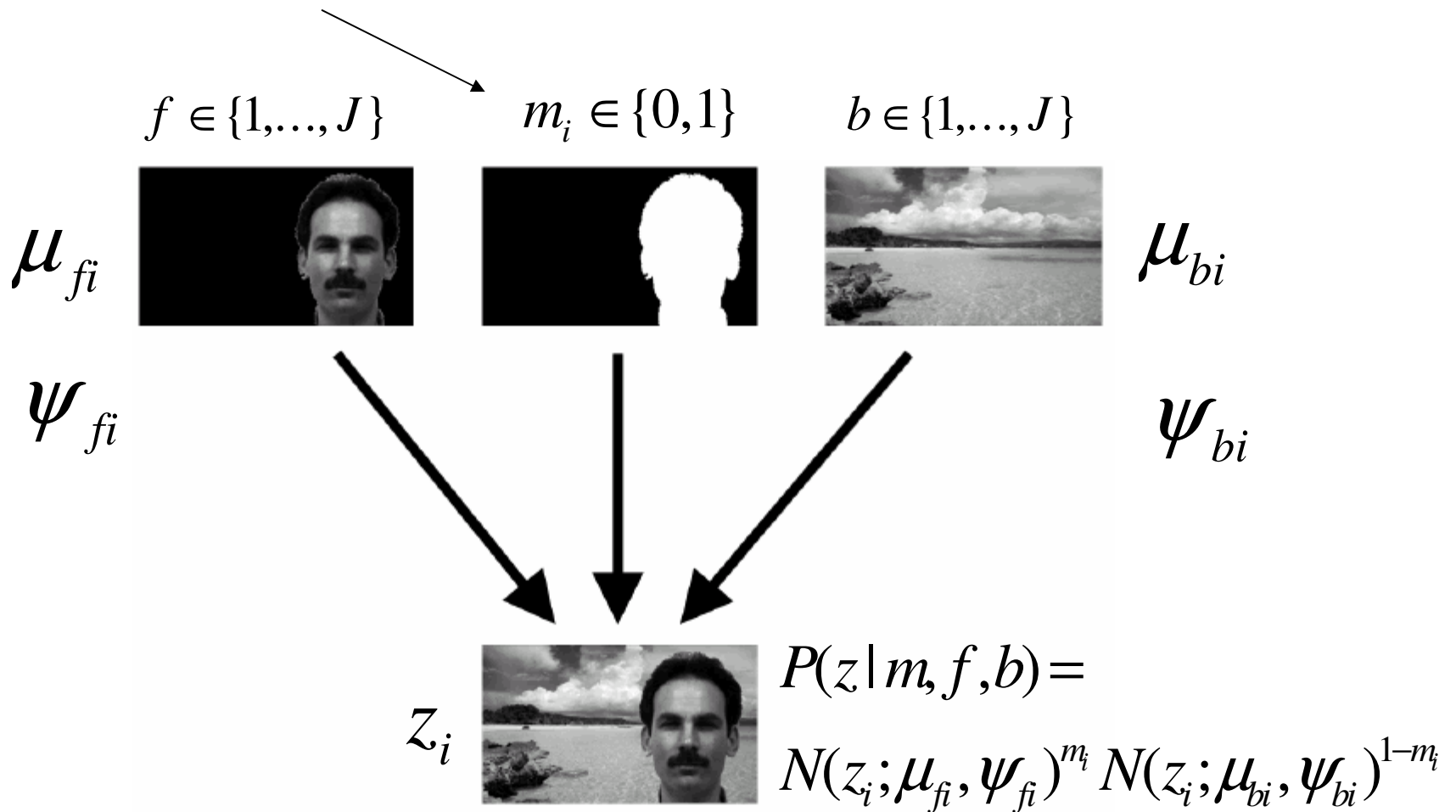
# Generative model



Let  $\alpha_i$  be the probability that  $m_i=1$  given that the foreground class is  $f$ ,  
 i.e.,  $P(m_i=1 | f) = \alpha_i$ ,  $P(m_i=0 | f) = 1-\alpha_i$ ,

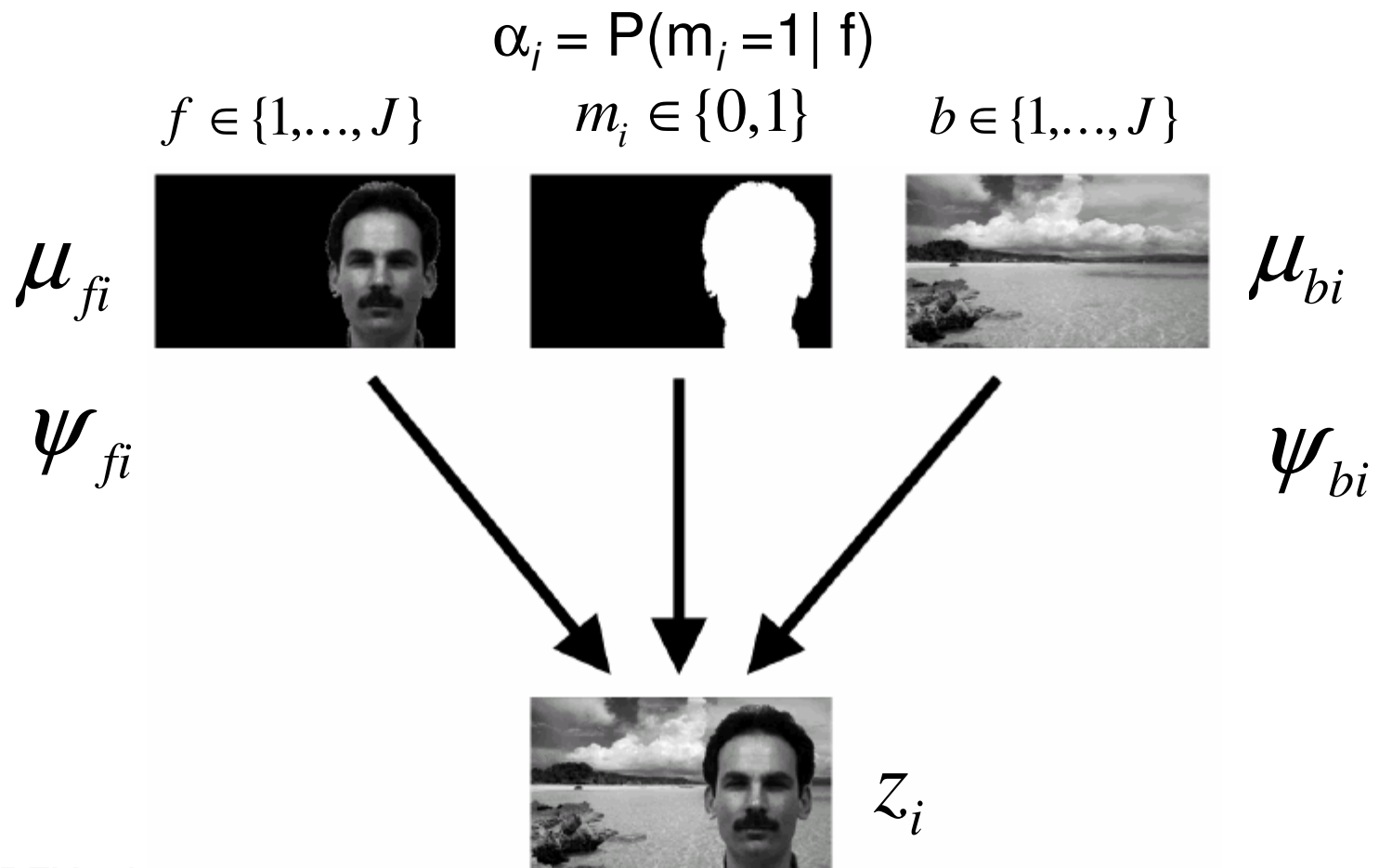


Let  $\alpha_i$  be the probability that  $m_i=1$  given that the foreground class is  $f$ ,  
 i.e.,  $P(m_i=1 | f) = \alpha_i$ ,  $P(m_i=0 | f) = 1-\alpha_i$ ,



The joint probability distribution

$$P(z, m, f, b) = P(b)P(f)\left(\prod_{i=1}^K P(m_i | f)\right)\left(\prod_{i=1}^K P(z_i | m_i, f, b)\right).$$



The joint probability distribution

$$P(z, m, f, b) = P(b)P(f)\left(\prod_{i=1}^K P(m_i | f)\right)\left(\prod_{i=1}^K P(z_i | m_i, f, b)\right)$$

Because  $m$  is binary, we can write:

$$P(z_i | m_i, f, b) = P^f(z_i | f)^{m_i} P^b(z_i | b)^{1-m_i}$$

$$P(z, m, f, b) = \pi_b \pi_f \left( \prod_{i=1}^K \alpha_{fi}^{m_i} (1 - \alpha_{fi})^{1-m_i} N(z_i; \mu_{fi}, \psi_{fi})^{m_i} N(z_i; \mu_{bi}, \psi_{bi})^{1-m_i} \right).$$

Mask at pixel  $i$       Observed intensity of the pixel  $i$

Prior distribution of foreground and background classes

Probability of occlusion

Means and variances of the layer classes

# Joint distribution over all variables and parameters in the dataset

- We assume uniform prior over the parameters

$$P(\mu, \psi, \pi, \alpha, f^{(1)}, b^{(1)}, m^{(1)}, \dots, f^{(T)}, b^{(T)}, m^{(T)} | z^{(1)}, \dots, z^{(T)}) \\ \propto \prod_{t=1}^T \left( \pi_{f^{(t)}} \pi_{b^{(t)}} \left( \prod_{i=1}^K \alpha_{f^{(t)}i}^{m_i^{(t)}} (1 - \alpha_{f^{(t)}i})^{1-m_i^{(t)}} \mathcal{N}(z_i^{(t)}; \mu_{f^{(t)}i}, \psi_{f^{(t)}i})^{m_i^{(t)}} \mathcal{N}(z_i^{(t)}; \mu_{b^{(t)}i}, \psi_{b^{(t)}i})^{1-m_i^{(t)}} \right) \right)$$

# Inference

$$P(m, f, b | z) \propto \pi_b \pi_f \left( \prod_{i=1}^K \alpha_{f_i}^{m_i} (1 - \alpha_{f_i})^{1-m_i} N(z_i; \mu_{f_i}, \psi_{f_i})^{m_i} N(z_i; \mu_{b_i}, \psi_{b_i})^{1-m_i} \right).$$

$$P(m, f, b | z) = P(f, b | z) P(m | f, b, z) = P(f, b | z) \prod_{i=1}^K P(m_i | f, b, z).$$

$$\begin{aligned} P(f, b | z) &\propto P(f, b, z) = \sum_{m_1} \cdots \sum_{m_K} P(m, f, b, z) \\ &= \pi_b \pi_f \prod_{i=1}^K (\alpha_{f_i} N(z_i; \mu_{f_i}, \psi_{f_i}) + (1 - \alpha_{f_i}) N(z_i; \mu_{b_i}, \psi_{b_i})). \end{aligned}$$



# Bayesian Networks

B. J. FREY & N. JOJIC

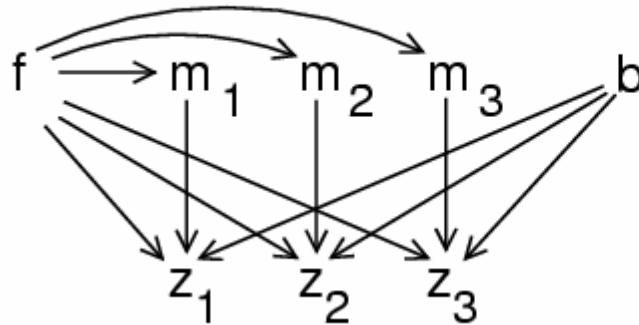
# Bayesian network

- *MAY* be constructed using knowledge of causal relationships
- Quickly conveys the factorization of a distribution
- Clearly expresses dependencies and independencies between variables
- Can be used to derive fast inference algorithms

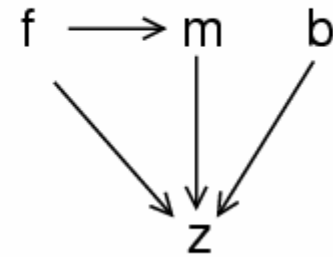
# Bayes net for occlusion model

$$P(z, m, f, b) = P(b)P(f)\left(\prod_{i=1}^K P(m_i | f)\right)\left(\prod_{i=1}^K P(z_i | m_i, f, b)\right)$$

(a)

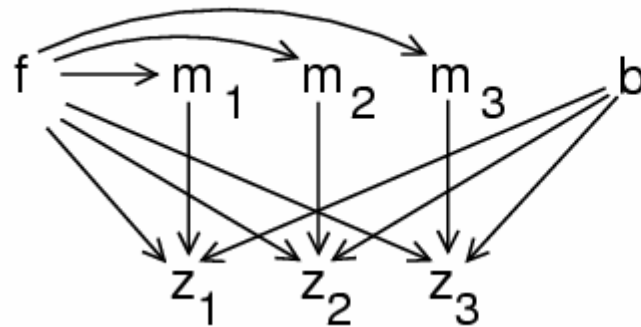


(b)

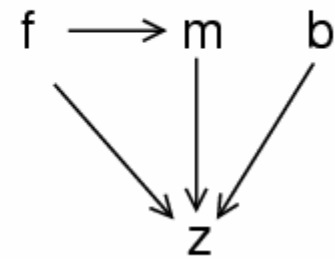


# Simulating Bayes nets

(a)

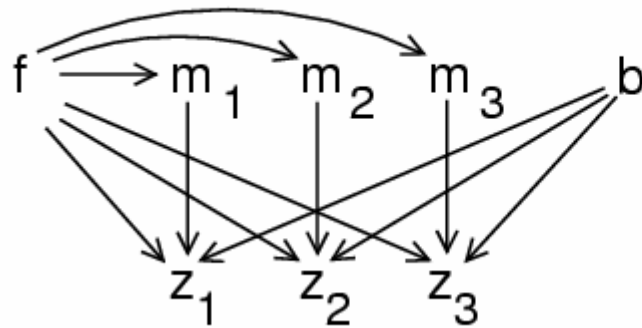


(b)

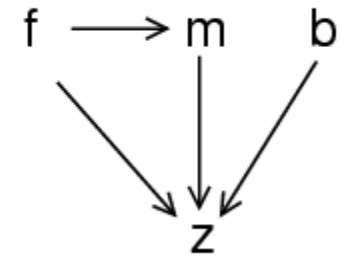


# Conditional independence

(a)



(b)



# A greedy inference and learning algorithm: ICM

# Iterative Conditional Modes (ICM)

- For  $t = 1, \dots, T$

$$\{ f^{(t)} \leftarrow \operatorname{argmax}_{f^{(t)}} [\pi_{f^{(t)}} \prod_{i:m_i^{(t)}=1} \mathcal{N}(z_i^{(t)}; \mu_{f^{(t)}i}, \psi_{f^{(t)}i})]$$

$$\{ b^{(t)} \leftarrow \operatorname{argmax}_{b^{(t)}} [\pi_{b^{(t)}} \prod_{i:m_i^{(t)}=0} \mathcal{N}(z_i^{(t)}; \mu_{b^{(t)}i}, \psi_{b^{(t)}i})]$$

$$\{ \text{For } i = 1, \dots, K: m_i^{(t)} \leftarrow \begin{cases} 1 & \text{if } \alpha_{f^{(t)}i} \mathcal{N}(z_i^{(t)}; \mu_{f^{(t)}i}, \psi_{f^{(t)}i}) > (1 - \alpha_{f^{(t)}i}) \mathcal{N}(z_i^{(t)}; \mu_{b^{(t)}i}, \psi_{b^{(t)}i}) \\ 0 & \text{otherwise} \end{cases}$$

- For  $j = 1, \dots, J$

$$\{ \pi_j \leftarrow (\sum_{t=1}^T [f^{(t)} = j] + \sum_{t=1}^T [b^{(t)} = j]) / 2T$$

- For  $j = 1, \dots, J$ , for  $i = 1, \dots, K$

$$\{ \alpha_{ji} \leftarrow (\sum_{t=1}^T [f^{(t)} = j] m_i^{(t)}) / (\sum_{t=1}^T [f^{(t)} = j])$$

$$\{ \mu_{ji} \leftarrow (\sum_{t=1}^T [f^{(t)} = j \text{ or } b^{(t)} = j] z_i^{(t)}) / (\sum_{t=1}^T [f^{(t)} = j \text{ or } b^{(t)} = j])$$

$$\{ \psi_{ji} \leftarrow (\sum_{t=1}^T [f^{(t)} = j \text{ or } b^{(t)} = j] (z_i^{(t)} - \mu_{ji})^2) / (\sum_{t=1}^T [f^{(t)} = j \text{ or } b^{(t)} = j])$$

Here, the Iverson notation is used where [True] = 1 and [False] = 0.

# A sample from the dataset

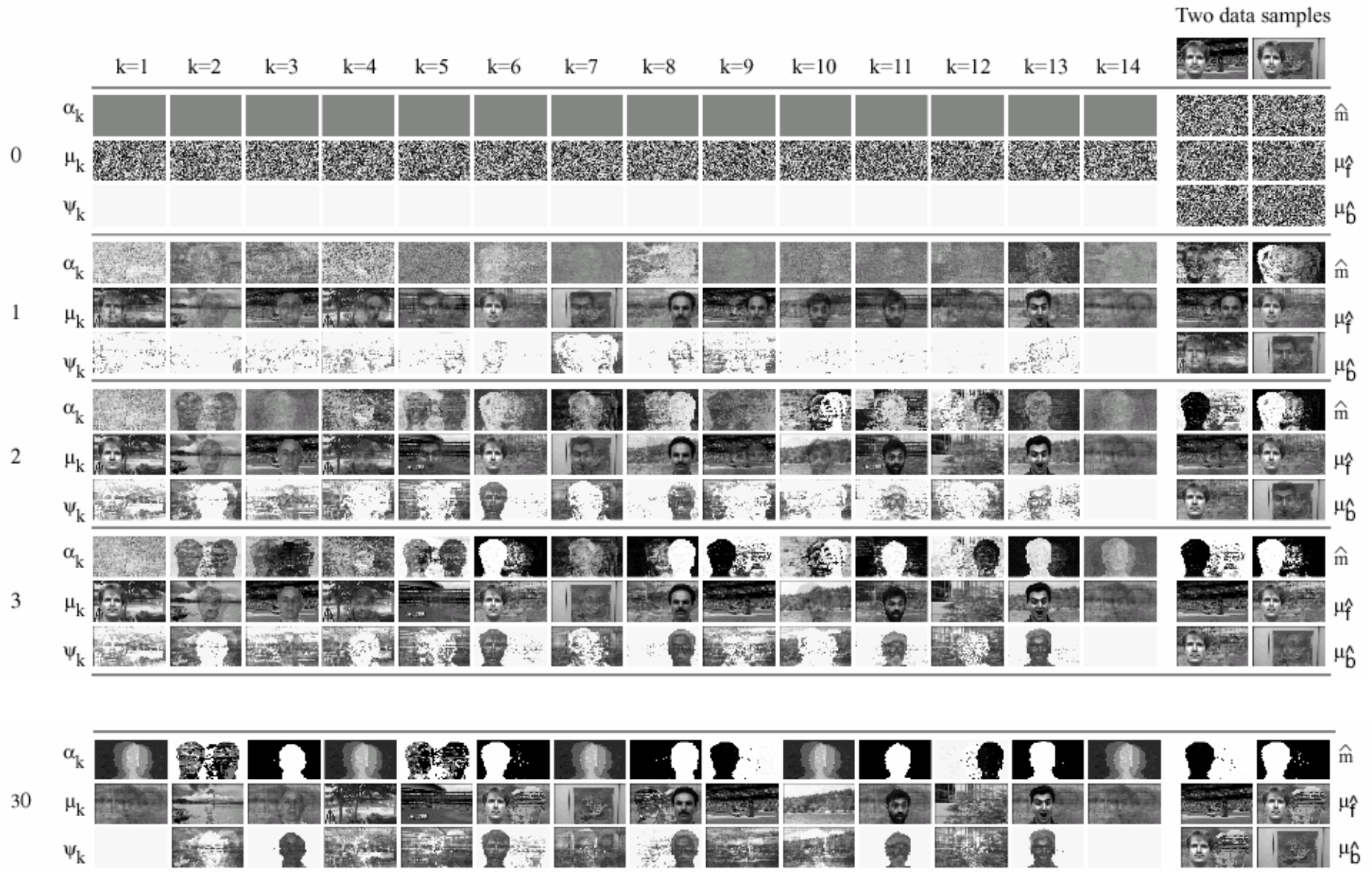
(5 different faces against 7 different backgrounds)



B. J. FREY & N. JOJIC

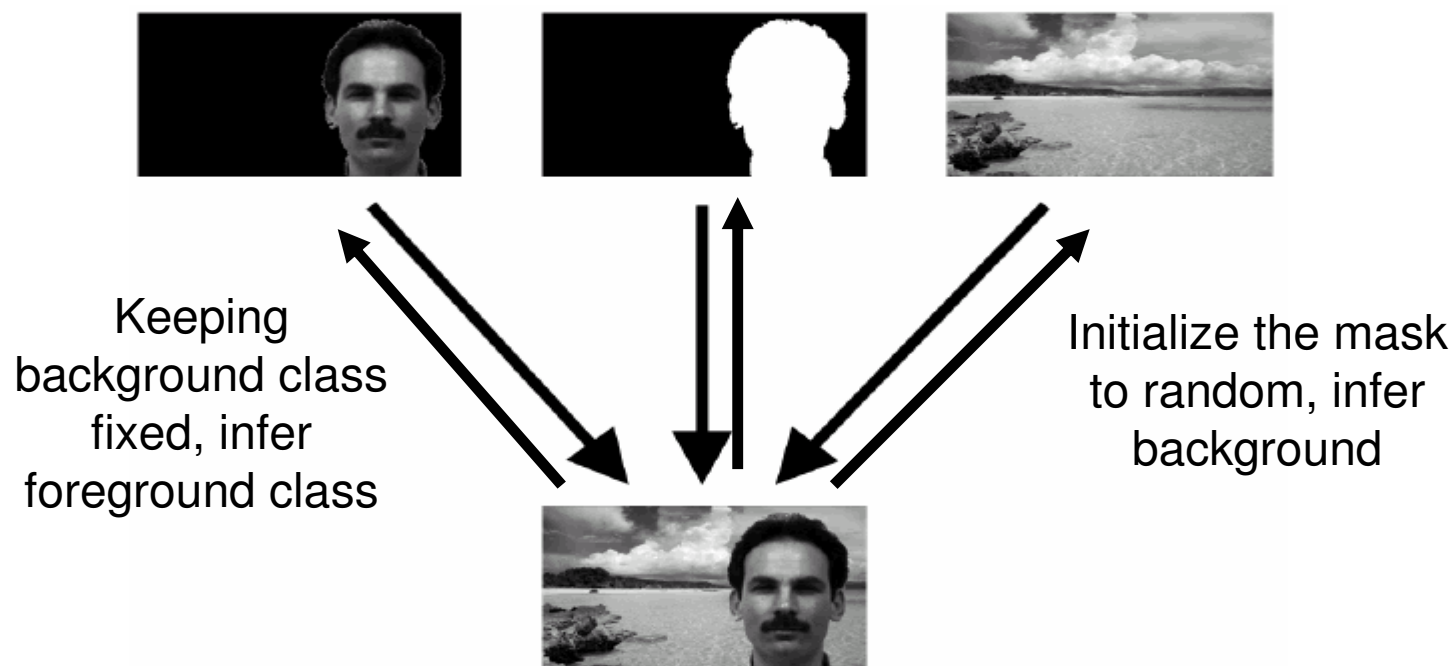


# Iterative model optimization



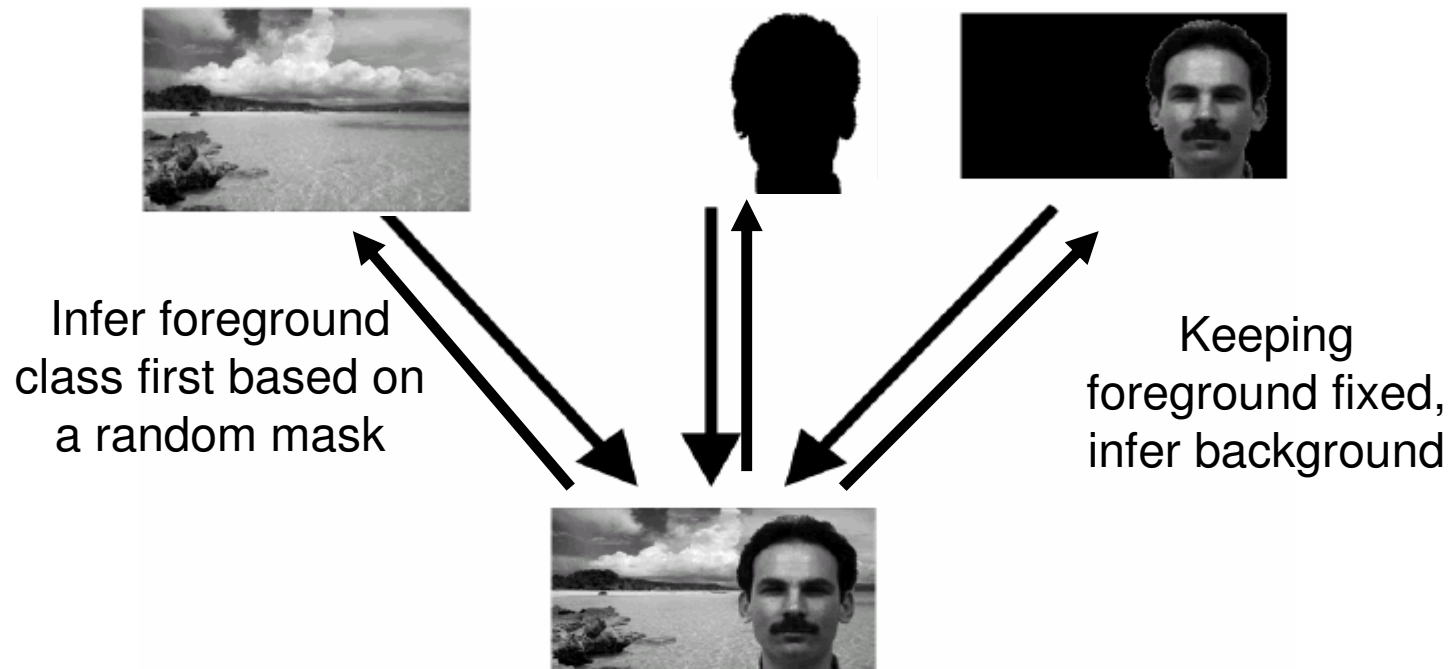
# Inferring hidden variables using ICM

- Given the class means and variances and mask probabilities for each class, for the shown dataset, it is (sometimes) possible to invert the generative process













































# However...

- Order of updates is important



# Inferring the hidden *and* estimating the parameters of a 14-class model using ICM

Class j	Class parameters		
	$\alpha$	$\mu$	$\psi$
1			
2			
3			
4			
5			
6			
7			
8			
9			
10			
11			
12			
13			
14			

# ICM

- Simple, fast
- Prone to local maxima
- Can be seen as iteratively employing simplified conditional posteriors, e.g.,

$$P(m, f, b | z) \propto \pi_b \pi_f \prod_{i=1}^K \alpha_{f_i}^{m_i} (1 - \alpha_{f_i})^{1 - m_i} N(z_i; \mu_{f_i}, \psi_{f_i})^{m_i} N(z_i; \mu_{b_i}, \psi_{b_i})^{1 - m_i}$$

$$P(b | \hat{f}, \hat{m}, z) \propto \pi_b \pi_{\hat{f}} \prod_{i=1}^K \alpha_{\hat{f}_i}^{\hat{m}_i} (1 - \alpha_{\hat{f}_i})^{1 - \hat{m}_i} N(z_i; \mu_{\hat{f}_i}, \psi_{\hat{f}_i})^{\hat{m}_i} N(z_i; \mu_{b_i}, \psi_{b_i})^{1 - \hat{m}_i}$$

$$\hat{b} = \arg \max P(b | \hat{f}, \hat{m}, z)$$

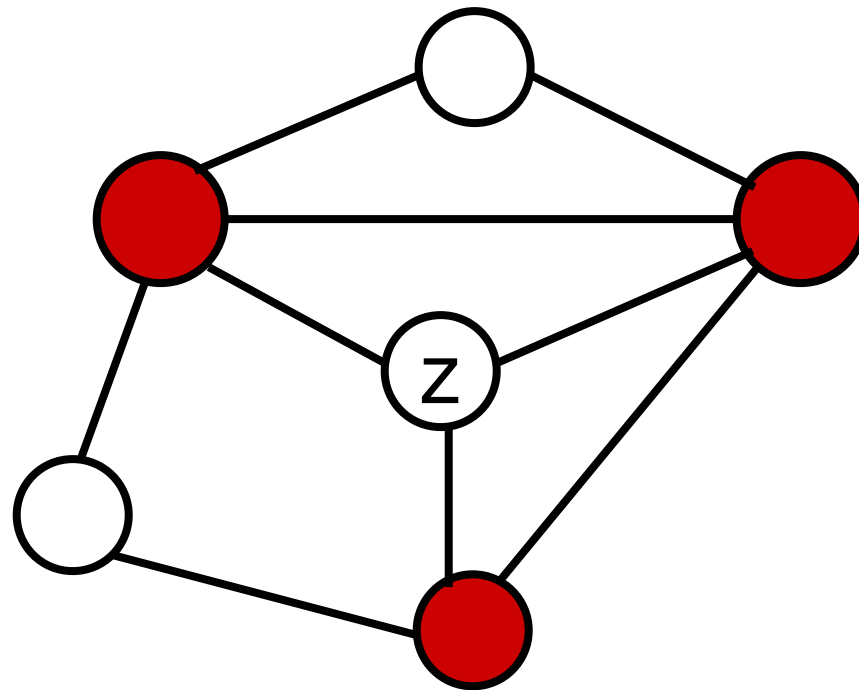
- Related to variational methods that use less severe approximations of the posterior

# Markov Random Fields (MRFs)

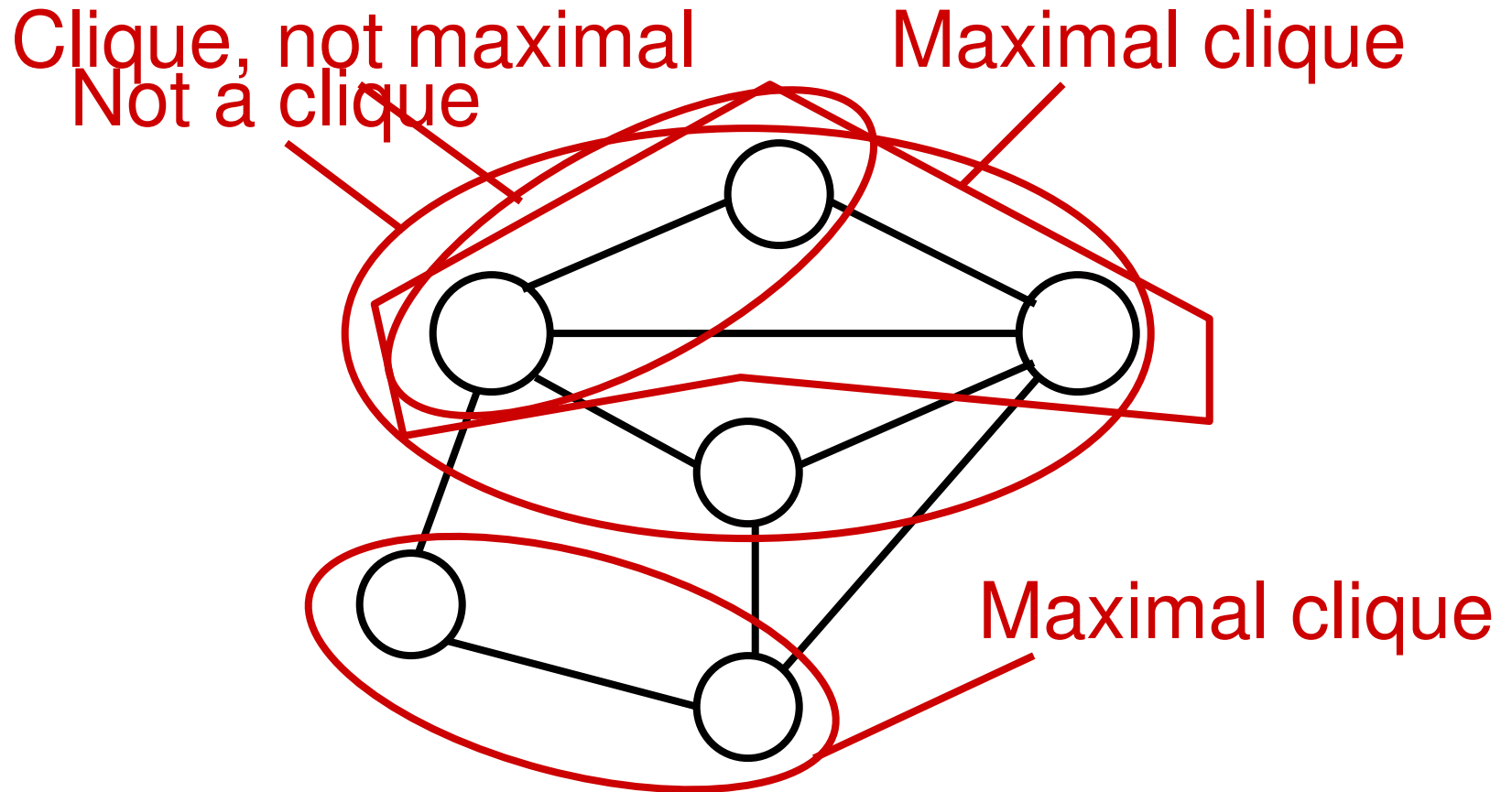
B. J. FREY & N. JOJIC

# Markov random fields (MRFs)

- Undirected graph on variables
- Each variable is independent of all other variables, given its neighbors



# Cliques and maximal cliques





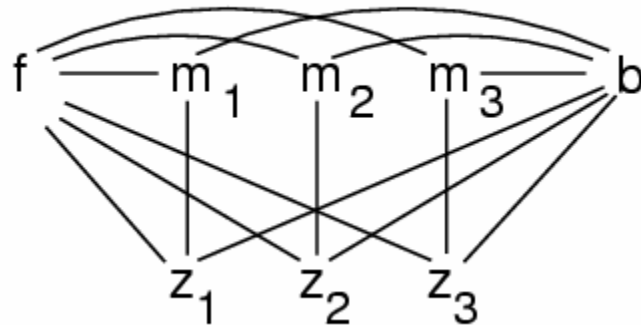
# The distribution for an MRF

$$P(x_1, \dots, x_N) = \alpha \prod_i \Psi_i(x_{C_i})$$

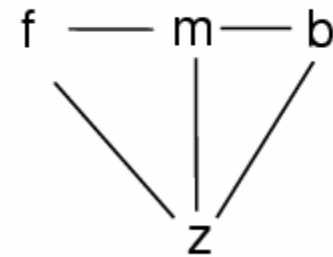
- $i$  = index of maximal clique
  - $C_i$  = set of indices of variables
  - $x_{C_i}$  = set of variables
- $\Psi_i(x_{C_i})$  is a “potential” (“local function”)
- $\alpha$  is a normalizing constant

# Occlusion model

(c)



(d)



$$P(z, m, f, b) = P(b)P(f)\left(\prod_{i=1}^K P(m_i | f)\right)\left(\prod_{i=1}^K P(z_i | m_i, f, b)\right).$$

# Factor Graphs

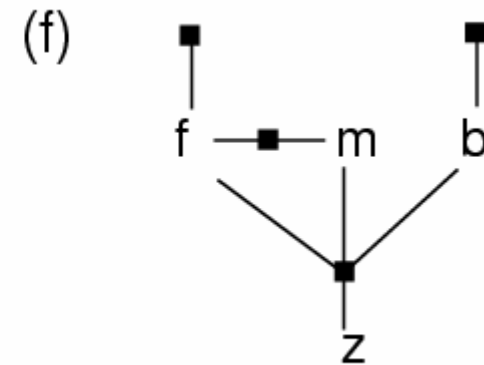
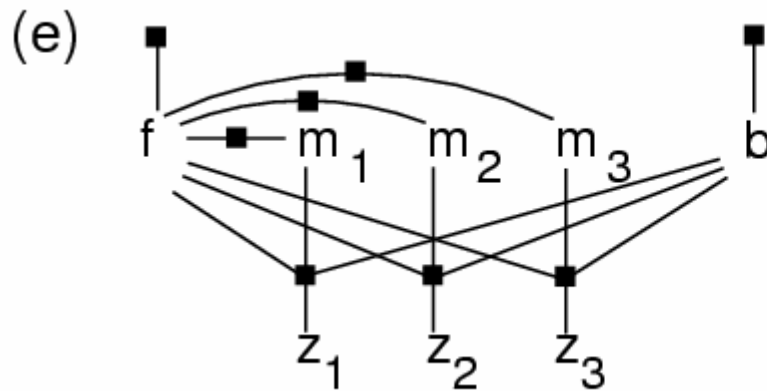
B. J. FREY & N. JOJIC

# Factor graphs

(Kschischang, Frey and Loeliger, 1999)

- Bipartite graph: nodes for variables and functions
- A *potential (local function)* is associated with each function node – this function depends on the neighboring variables
- The joint distribution is given by the product of the local functions
- Edges can be directed or undirected (a directed edge indicates a conditional probability)

# Occlusion model



$$P(z, m, f, b) = P(b)P(f)\left(\prod_{i=1}^K P(m_i | f)\right)\left(\prod_{i=1}^K P(z_i | m_i, f, b)\right).$$

# Parameterized Models and Bayesian Learning

B. J. FREY & N. JOJIC

# Parameters as variables

- Recall for our toy problem:

$$P(z, m, f, b) = \pi_b \pi_f \left( \prod_{i=1}^K \alpha_{fi}^{m_i} (1 - \alpha_{fi})^{1-m_i} N(z_i; \mu_{fi}, \psi_{fi})^{m_i} N(z_i; \mu_{bi}, \psi_{bi})^{1-m_i} \right).$$

- Interpret the parameters as variables, where the joint distribution is

$$P(z, m, f, b, \pi, \alpha, \mu, \psi) = P(b | \pi) P(f | \pi) P(m | f, \alpha) P(z | m, f, b, \mu, \psi) P(\pi) P(\alpha) P(\mu) P(\psi).$$

where

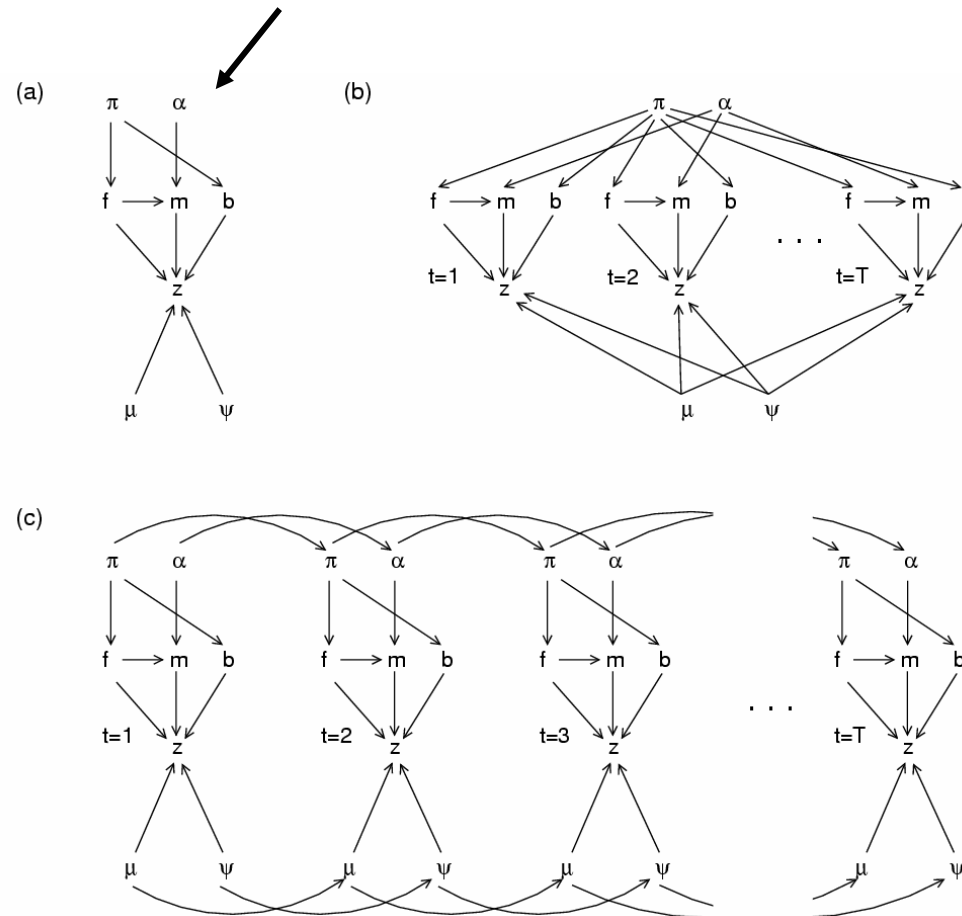
$$P(b | \pi) = \pi_b \quad P(f | \pi) = \pi_f \quad P(m_i | f, \alpha_{1i}, \dots, \alpha_{Ji}) = \alpha_{fi}^{m_i} (1 - \alpha_{fi})^{1-m_i}$$

$$P^f(z_i | f, \mu_{1i}, \dots, \mu_{Ji}, \psi_{1i}, \dots, \psi_{Ji}) = N(z_i; \mu_{fi}, \psi_{fi})$$

$$P^b(z_i | b, \mu_{1i}, \dots, \mu_{Ji}, \psi_{1i}, \dots, \psi_{Ji}) = N(z_i; \mu_{bi}, \psi_{bi})$$

# Bayes nets for Bayesian version of toy problem

$$P(z, m, f, b, \pi, \alpha, \mu, \psi) = P(b | \pi)P(f | \pi)P(m | f, \alpha)P(z | m, f, b, \mu, \psi)P(\pi)P(\alpha)P(\mu)P(\psi).$$





# Introducing training data

- Visible variables:  $v = (v^{(1)}, \dots, v^{(T)})$ 
  - $(t)$  denotes the  $t^{\text{th}}$  training case
- Hidden variables:  $h = (h^\theta, h^{(1)}, \dots, h^{(T)})$ 
  - Parameters:  $h^\theta$
  - Variables for the  $t^{\text{th}}$  training case:  $h^{(t)}$
- Joint distribution:  $P(h, v) = P(h^\theta) \prod_{t=1}^T P(h^{(t)}, v^{(t)} | h^\theta)$ .
- Parameter prior:  $P(h^\theta)$
- Parameter likelihood:  $\prod_{t=1}^T P(h^{(t)}, v^{(t)} | h^\theta)$

# Parameter priors

- Uniform priors
  - Computationally attractive
  - Problem: “Uniform” is inconsistent w.r.t. reparameterization
- Conjugate priors
  - Prior x Likelihood has same form as likelihood
  - Conjugate prior can be thought of as likelihood for “fake” data, eg, prior counts in a coin toss experiment

# Algorithms for Inference and Learning

B. J. FREY & N. JOJIC

# Problem set-up

- Probabilistic inference and learning entail computing the intractable distribution,

$$P(h | v) = \frac{P(h, v)}{\int_h P(h, v)},$$

- Note that w.r.t.  $h$ ,  $P(h | v) \propto P(h, v)$ .
- In a graphical model,  $P(h, v)$  factorizes

# General “brute force” inference

- Suppose  $x_1, x_2, \dots, x_N$  are binary

$$P(x_1) = \sum_{x_2} \sum_{x_3} \dots \sum_{x_N} P(x_1, x_2, \dots, x_N)$$

- This takes about  $2^N$  operations
- Generally, computing  $P(x_i | \text{Observed } x\text{'s})$  takes  $2^{(N - \text{\#observed } x\text{'s})}$  operations

# Exact posterior for occlusion model

$$P(\mu, \psi, \pi, \alpha, f^{(1)}, b^{(1)}, m^{(1)}, \dots, f^{(T)}, b^{(T)}, m^{(T)} | z^{(1)}, \dots, z^{(T)})$$
$$\propto \prod_{t=1}^T \left( \pi_{f^{(t)}} \pi_{b^{(t)}} \left( \prod_{i=1}^K \alpha_{f^{(t)}i}^{m_i^{(t)}} (1 - \alpha_{f^{(t)}i})^{1-m_i^{(t)}} \mathcal{N}(z_i^{(t)}; \mu_{f^{(t)}i}, \psi_{f^{(t)}i})^{m_i^{(t)}} \mathcal{N}(z_i^{(t)}; \mu_{b^{(t)}i}, \psi_{b^{(t)}i})^{1-m_i^{(t)}} \right) \right)$$

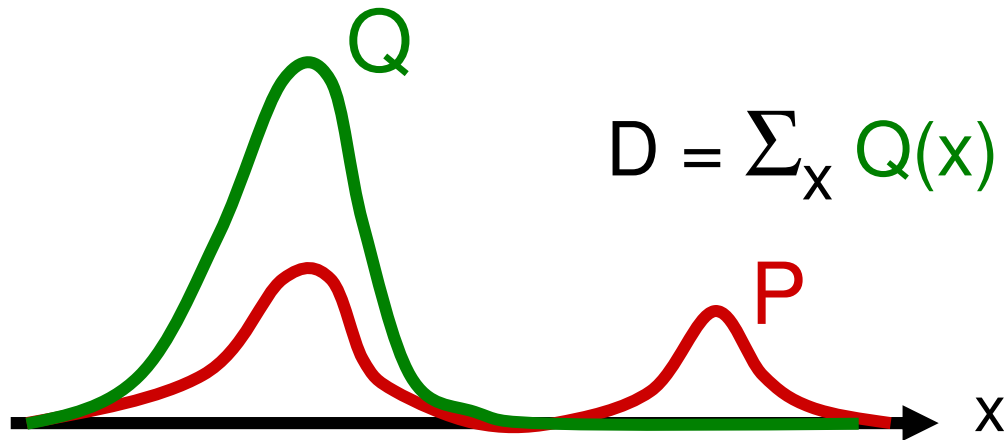
# Approximate inference

- Since  $P(h|v)$  is intractable, we search for a surrogate,  $Q(h)$  that is tractable
- Measure of quality of  $Q$ : Kullback-Leibler divergence between  $Q$  and  $P$ ,

$$D(Q, P) = \int_h Q(h) \log \frac{Q(h)}{P(h|v)}.$$

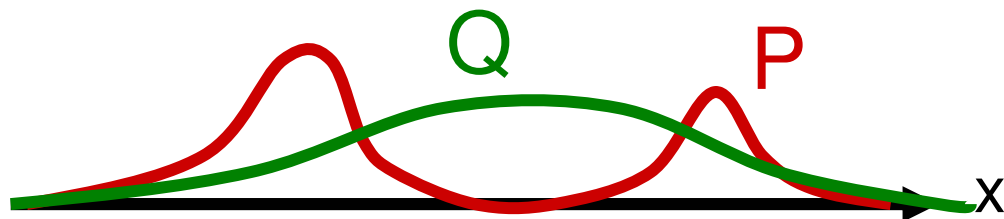
- Properties:
  - $D \geq 0$
  - $D = 0$  iff  $Q(h) = P(h|v)$

## More on the “distance”



$$D = \sum_x Q(x) \log[Q(x)/P(x)]$$

For unimodal Q, minimizing D favors this Q



$$D = \sum_x P(x) \log[P(x)/Q(x)]$$

For unimodal Q, minimizing D favors this Q



# Free energy

- In fact,  $D$  is intractable
- We can make it tractable by subtracting  $\log P(v)$

$$F(Q, P) = D(Q, P) - \log P(v)$$

$$= \int_h Q(h) \log \frac{Q(h)}{P(h|v)} - \int_h Q(h) \log P(v)$$

$$= \int_h Q(h) \log \frac{Q(h)}{P(h|v)P(v)}$$

$$= \int_h Q(h) \log \frac{Q(h)}{P(h, v)}$$

Factorizes, for  
a graphical model



# Alternative derivation of free energy

- $\log P(v)$  is the log-likelihood of the data
- Computing  $\log P(v)$  is intractable, but we can bound it using Jensen's inequality:

$$\begin{aligned}\log P(v) &= \log\left(\int_h P(h, v)\right) \\ &= \log\left(\int_h Q(h) \frac{P(h, v)}{Q(h)}\right) \\ &\geq \int_h Q(h) \log\left(\frac{P(h, v)}{Q(h)}\right) = -F(Q, P).\end{aligned}$$

# Properties of free energy

- $F \geq -\log P(v)$
- The minimum of  $F$  w.r.t  $Q$  gives

$$F = -\log P(v)$$

$$Q(h) = P(h|v)$$

# Approaches to constructing $Q$

- Parameterize  $Q$  directly
  - Iterative conditional modes
  - Variational methods (mean field)
  - Structured variational methods
  - Hybrids (the EM algorithm)
- Parameterize marginals of  $Q$ 
  - Bethe approximation (probability propagation)
  - Kikuchi approximation (generalized prob prop)
- Represent  $Q$  using samples
  - Monte Carlo
  - Markov chain Monte Carlo

# I.I.D. Training cases

## Free energy for i.i.d. training cases

From (5), for a training set of  $T$  i.i.d. training cases with hidden variables  $h = (h^\theta, h^{(1)}, \dots, h^{(T)})$  and visible variables  $v = (v^{(1)}, \dots, v^{(T)})$ , we have  $P(h, v) = P(h^\theta) \prod_{t=1}^T P(h^{(t)}, v^{(t)} | h^\theta)$ . The free energy is

$$\begin{aligned} F(Q, P) &= \int_h Q(h) \log Q(h) - \int_h Q(h) \log P(h, v) \\ &= \int_h Q(h) \log Q(h) - \int_{h^\theta} Q(h^\theta) \log P(h^\theta) - \sum_{t=1}^T \int_{h^{(t)}, h^\theta} Q(h^{(t)}, h^\theta) \log P(h^{(t)}, v^{(t)} | h^\theta). \end{aligned} \quad (8)$$

The decomposition of  $F$  into a sum of one term for each training case simplifies learning.

# Point inference

- Discrete hidden variables:  $\hat{h}$  is a point estimate

$$F(Q, P) = \sum_h [h = \hat{h}] \log \frac{[h = \hat{h}]}{P(h, v)} = -\log P(\hat{h}, v).$$

- Minimizing  $F$  w.r.t.  $\hat{h}$  Maximizes  $P(\hat{h}, v)$  w.r.t.  $h$

- Continuous hidden variables:

$$F(Q, P) = \int_h \delta(h - \hat{h}) \log \frac{\delta(h - \hat{h})}{P(h, v)} = -\log P(\hat{h}, v) - H_\delta,$$

- Very common in engineering and science
- Problem:  $H_\delta \rightarrow -\infty$

# Iterative Conditional Modes (ICM)

**Initialization.** Pick values for all hidden variables  $h$  (randomly, or cleverly).

**ICM Step.** Select one of the hidden variables,  $h_i$ . Holding all other variables constant, set  $h_i$  to its MAP value:

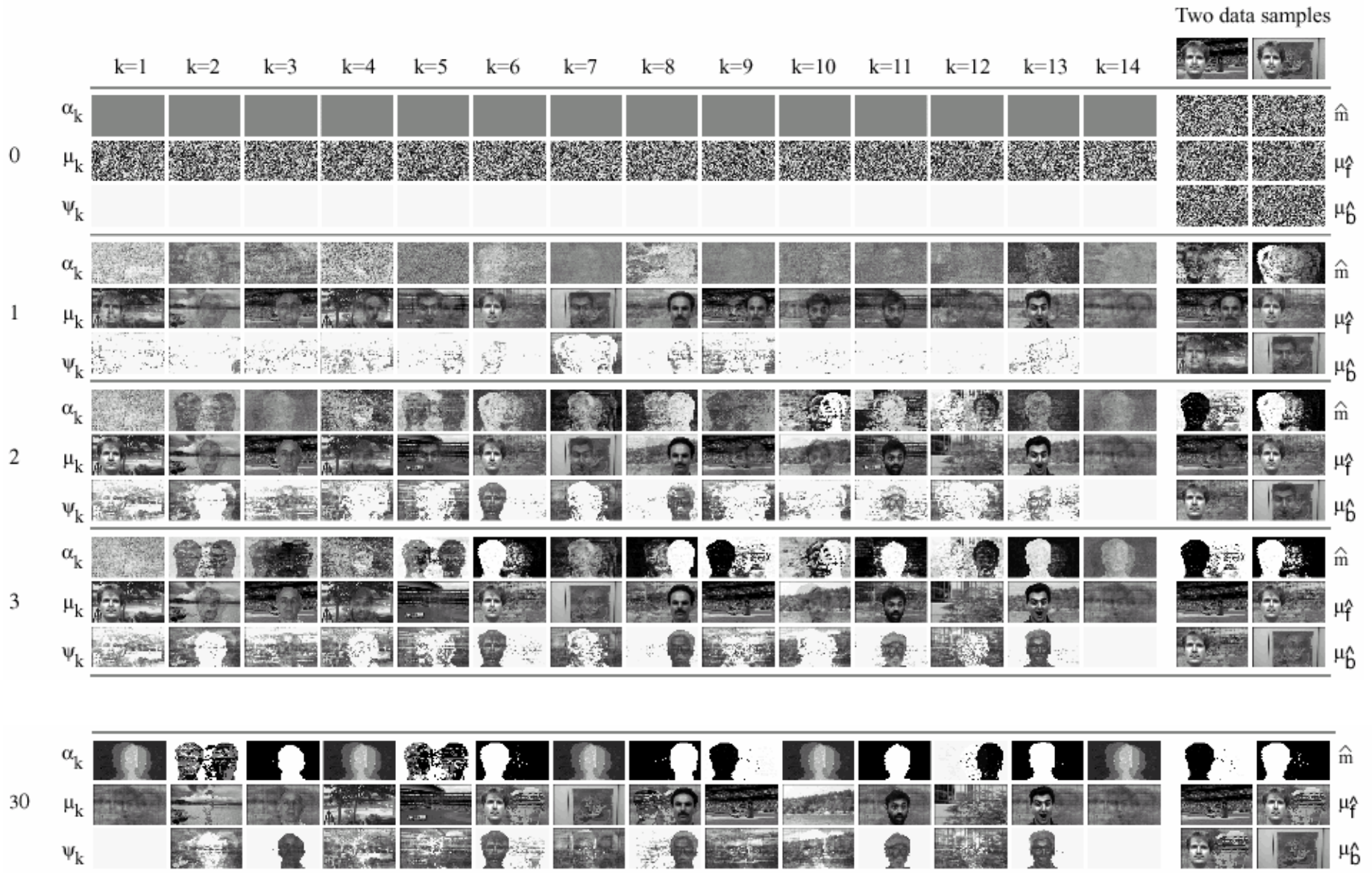
$$h_i \leftarrow \operatorname{argmax}_{h_i} P(h_i | h \setminus h_i, v) = \operatorname{argmax}_{h_i} P(h, v),$$

where  $h \setminus h_i$  is the set of all hidden variables other than  $h_i$ .

**Repeat for a fixed number of iterations or until convergence.**

- Examples
  - ICM in mixture of Gaussians = k-means clustering
  - Clever ICM in M of G = sequential k-means clustering

# Recall ICM for occlusion model





# Properties of ICM

- Fast, greedy
- Can get stuck in local minima
- Does not account for mass of mode
- Does not account for multiple modes
- Speed depends on choice of substructures
- Usually easy to implement

# The EM Algorithm

- Use the following Q-distribution:

$$Q(h) = \delta(h^\theta - \hat{h}^\theta) Q(h^{(1)}, \dots, h^{(T)}).$$

- IID data: The training cases are independent, so

$$Q(h) = \delta(h^\theta - \hat{h}^\theta) \prod_{t=1}^T Q(h^{(t)}).$$

- Free energy:

$$F(Q, P) = -\log P(\hat{h}^\theta) + \sum_{t=1}^T \left( \int_{h^{(t)}} Q(h^{(t)}) \log \frac{Q(h^{(t)})}{P(h^{(t)}, v^{(t)} | \hat{h}^\theta)} \right).$$

# The EM Algorithm

**Initialization.** Choose values for the parameters,  $\hat{h}^\theta$  (randomly, or cleverly).

**E Step.** Minimize  $F(Q, P)$  w.r.t.  $Q$  by setting

$$Q(h^{(t)}) \leftarrow P(h^{(t)}|v^{(t)}, \hat{h}^\theta),$$

for each training case, given the parameters  $\hat{h}^\theta$  and the data  $v^{(t)}$ .

**M Step.** Minimize  $F(Q, P)$  w.r.t. the model parameters  $\hat{h}^\theta$  by solving

$$-\frac{\partial}{\partial \hat{h}^\theta} \log P(\hat{h}^\theta) - \sum_{t=1}^T \left( \int_{h^{(t)}} Q(h^{(t)}) \frac{\partial}{\partial \hat{h}^\theta} \log P(h^{(t)}, v^{(t)}|\hat{h}^\theta) \right) = 0. \quad (9)$$

For  $M$  parameters, this is a system of  $M$  equations. Often, the prior on the parameters is assumed to be uniform,  $P(\hat{h}^\theta) = \text{const}$ , in which case the first term in the above expression vanishes.

**Repeat for a fixed number of iterations or until convergence.**

Occlusion model:  
EM as an algorithm for minimizing the free energy

$$\begin{aligned}
 F(Q, P) &= \int_h Q(h) \log Q(h) - \int_h Q(h) \log P(h, v) \\
 &= \int_h Q(h) \log Q(h) - \int_{h^\theta} Q(h^\theta) \log P(h^\theta) - \\
 &\quad - \sum_{t=1}^T \int_{h^{(t)}, h^\theta} Q(h^{(t)}, h^\theta) \log P(h^{(t)}, v^{(t)} | h^\theta).
 \end{aligned}$$

$$P \propto \prod_{t=1}^T (\pi_{f^{(t)}} \pi_{b^{(t)}} \left( \prod_{i=1}^K \alpha_{f^{(t)}i}^{m_i^{(t)}} (1 - \alpha_{f^{(t)}i})^{1-m_i^{(t)}} N(z_i^{(t)}; \mu_{f^{(t)}i}, \psi_{f^{(t)}i})^{m_i^{(t)}} N(z_i^{(t)}; \mu_{b^{(t)}i}, \psi_{b^{(t)}i})^{1-m_i^{(t)}} \right)).$$

$$\begin{aligned}
 Q &= \left( \prod_k \delta(\pi_k - \hat{\pi}_k) \right) \left( \prod_{k,i} \delta(\mu_{ki} - \hat{\mu}_{ki}) \right) \left( \prod_{k,i} \delta(\psi_{ki} - \hat{\psi}_{ki}) \right) \left( \prod_{k,i} \delta(\alpha_{ki} - \hat{\alpha}_{ki}) \right) \\
 &\quad \cdot Q(b, f) \prod_i Q(m_i | b, f).
 \end{aligned}$$

Occlusion model:  
EM as an algorithm for minimizing the free energy

**E**

$$Q(m_i = 1 | b, f) \leftarrow \frac{\alpha_{f_i} N(z_i; \mu_{f_i}, \psi_{f_i})}{\alpha_{f_i} N(z_i; \mu_{f_i}, \psi_{f_i}) + (1 - \alpha_{f_i}) N(z_i; \mu_{b_i}, \psi_{b_i})},$$

$$Q(b, f) \leftarrow c \pi_b \pi_f \exp\left\{-\sum_i (Q(m_i = 1 | b, f)) \left(\frac{(z_i - \mu_{f_i})^2}{2\psi_{f_i}} + \frac{\log 2\pi\psi_{f_i}}{2}\right) + (1 - Q(m_i = 1 | b, f)) \left(\frac{(z_i - \mu_{b_i})^2}{2\psi_{b_i}} + \frac{\log 2\pi\psi_{b_i}}{2}\right)\right\},$$

**M**

$$\pi_k \leftarrow (\sum_t Q(f^{(t)} = k) + \sum_t Q(b^{(t)} = k)) / (2T),$$

$$Q(b) \leftarrow \sum_f Q(b, f)$$

$$Q(f) \leftarrow \sum_b Q(b, f)$$

$$\alpha_{ki} \leftarrow \frac{\sum_t Q(m_i^{(t)} = 1, f^{(t)} = k)}{\sum_t Q(f^{(t)} = k)},$$

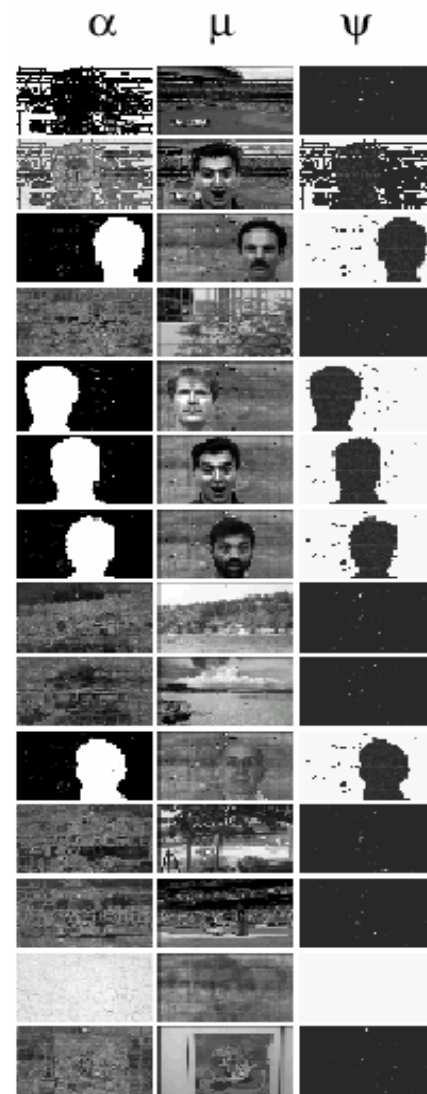
$$Q(m_i = 1, b) \leftarrow \sum_f Q(m_i = 1 | b, f) Q(b, f)$$

$$\mu_{ki} \leftarrow \frac{\sum_t (Q(m_i^{(t)} = 1, f^{(t)} = k) + Q(m_i^{(t)} = 0, b^{(t)} = k)) z_i^{(t)}}{\sum_t (Q(m_i^{(t)} = 1, f^{(t)} = k) + Q(m_i^{(t)} = 0, b^{(t)} = k))}.$$

$$Q(m_i = 1, f) \leftarrow \sum_b Q(m_i = 1 | b, f) Q(b, f)$$

$$\psi_{ki} \leftarrow \frac{\sum_t (Q(m_i^{(t)} = 1, f^{(t)} = k) + Q(m_i^{(t)} = 0, b^{(t)} = k)) (z_i^{(t)} - \mu_{ki})^2}{\sum_t (Q(m_i^{(t)} = 1, f^{(t)} = k) + Q(m_i^{(t)} = 0, b^{(t)} = k))}.$$

Toy problem:  
Result of EM learning



# Generalizations of EM

- Update  $h^\theta$  so as to only decrease F (not minimize F)
- Instead of updating Q using  $Q(h^{(t)}) \leftarrow P(h^{(t)} | v^{(t)}, \hat{h}^\theta)$  update Q so as to only decrease F
- Use another inference technique for the Q-distribution over hidden variables...

# Variational techniques

- Parameterize the Q-distribution:  $Q(h; \phi)$
- Now, the *variational* free energy is

$$F(Q, P) = \int_h Q(h; \phi) \log \frac{Q(h; \phi)}{P(h, v)}.$$

- Computational task: Min F w.r.t  $\phi$
- Choose form of Q so math works out
- Example: Mean Field,  $Q(h) = \prod_{i=1}^M Q(h_i),$



# Variational inference

**Initialization.** Pick values for the variational parameters,  $\phi$  (randomly, or cleverly).

**Optimization Step.** Decrease  $F(Q, P)$  by adjusting the parameter vector  $\phi$ , or a subset of  $\phi$ .

**Repeat for a fixed number of iterations or until convergence.**

# Variational EM

**Initialization.** Pick values for the variational parameters  $\phi^{(1)}, \dots, \phi^{(T)}$  and the model parameters  $\hat{h}^\theta$  (randomly, or cleverly).

**Generalized E Step.** Starting from the variational parameters from the previous iteration, modify  $\phi^{(1)}, \dots, \phi^{(T)}$  so as to decrease  $F$ .

**Generalized M Step.** Starting from the model parameters from the previous iteration, modify  $\hat{h}^\theta$  so as to decrease  $F$ .

**Repeat for a fixed number of iterations or until convergence.**

Occlusion model:  
Variational EM for minimizing the free energy

$$\begin{aligned}
 F(Q, P) &= \int_h Q(h) \log Q(h) - \int_h Q(h) \log P(h, v) \\
 &= \int_h Q(h) \log Q(h) - \int_{h^\theta} Q(h^\theta) \log P(h^\theta) - \\
 &\quad - \sum_{t=1}^T \int_{h^{(t)}, h^\theta} Q(h^{(t)}, h^\theta) \log P(h^{(t)}, v^{(t)} | h^\theta).
 \end{aligned}$$

$$P \propto \prod_{t=1}^T (\pi_{f^{(t)}} \pi_{b^{(t)}} (\prod_{i=1}^K \alpha_{f^{(t)}i}^{m_i^{(t)}} (1 - \alpha_{f^{(t)}i})^{1-m_i^{(t)}} N(z_i^{(t)}; \mu_{f^{(t)}i}, \psi_{f^{(t)}i})^{m_i^{(t)}} N(z_i^{(t)}; \mu_{b^{(t)}i}, \psi_{b^{(t)}i})^{1-m_i^{(t)}})).$$

$$\begin{aligned}
 Q &= (\prod_k \delta(\pi_k - \hat{\pi}_k)) (\prod_{k,i} \delta(\mu_{ki} - \hat{\mu}_{ki})) (\prod_{k,i} \delta(\psi_{ki} - \hat{\psi}_{ki})) (\prod_{k,i} \delta(\alpha_{ki} - \hat{\alpha}_{ki})) \\
 &\quad \cdot Q(b) Q(f) \prod_i Q(m_i). \qquad \text{VEM}
 \end{aligned}$$

Occlusion model:  
Variational EM for minimizing the free energy

$$\begin{aligned}
 F = & \sum_b Q(b) \log \frac{Q(b)}{\pi_b} + \sum_f Q(f) \log \frac{Q(f)}{\pi_f} \\
 & + \sum_i (q_i \log q_i + (1 - q_i) \log(1 - q_i)) - \sum_i (q_i (\sum_f Q(f) \log \alpha_{fi}) + (1 - q_i) (\sum_f Q(f) \log(1 - \alpha_{fi}))) \\
 & + \sum_i \sum_f Q(f) q_i \left( \frac{(z_i - \mu_{fi})^2}{2\psi_{fi}} + \frac{\log 2\pi\psi_{fi}}{2} \right) + \sum_i \sum_b Q(b) (1 - q_i) \left( \frac{(z_i - \mu_{bi})^2}{2\psi_{bi}} + \frac{\log 2\pi\psi_{bi}}{2} \right).
 \end{aligned}$$

$$P \propto \prod_{t=1}^T (\pi_{f^{(t)}} \pi_{b^{(t)}} \left( \prod_{i=1}^K \alpha_{f^{(t)}i}^{m_i^{(t)}} (1 - \alpha_{f^{(t)}i})^{1 - m_i^{(t)}} N(z_i^{(t)}; \mu_{f^{(t)}i}, \psi_{f^{(t)}i})^{m_i^{(t)}} N(z_i^{(t)}; \mu_{b^{(t)}i}, \psi_{b^{(t)}i})^{1 - m_i^{(t)}} \right)).$$

$$\begin{aligned}
 Q = & \left( \prod_k \delta(\pi_k - \hat{\pi}_k) \right) \left( \prod_{k,i} \delta(\mu_{ki} - \hat{\mu}_{ki}) \right) \left( \prod_{k,i} \delta(\psi_{ki} - \hat{\psi}_{ki}) \right) \left( \prod_{k,i} \delta(\alpha_{ki} - \hat{\alpha}_{ki}) \right) \\
 & \cdot Q(b) Q(f) \prod_i Q(m_i). \qquad \text{VEM}
 \end{aligned}$$

# Occlusion model: Comparison of exact EM with Variational EM

Class	$v_f$	$v_b$	$\alpha$	$\mu$	$\psi$	$v_f$	$v_b$	$\alpha$	$\mu$	$\psi$
1	0.04	0.11				0.06	0.07			
2	0	0.04				0.07	0.12			
3	0.20	0				0.07	0.04			
4	0	0.14				0.07	0.06			
5	0.19	0				0.07	0.15			
6	0.17	0				0.09	0.07			
7	0.19	0				0.10	0.03			
8	0	0.13				0.03	0.08			
9	0	0.13				0.09	0.06			
10	0.21	0				0.02	0.03			
11	0	0.12				0.10	0.04			
12	0	0.17				0.04	0.06			
13	0	0				0.07	0.15			
14	0	0.16				0.12	0.04			

Exact EM

Variational EM

Foreground and background frequencies

# Structured variational techniques

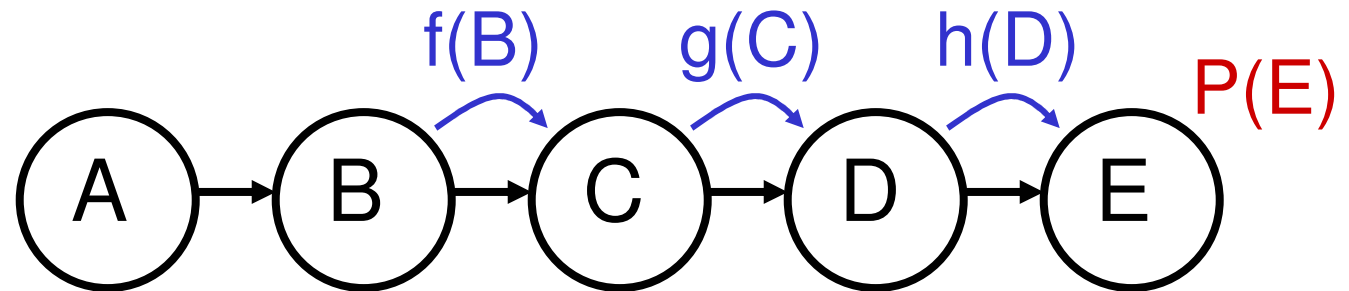
- We can also describe  $Q$  using a parameterized graphical model
- Choose the structure and form of  $Q$  so that

$$F(Q, P) = \int_h Q(h; \phi) \log \frac{Q(h; \phi)}{P(h, v)}.$$

is tractable

- Structured variational technique for toy problem: See tutorial paper

# The sum-product algorithm (belief propagation)



$$P(A,B,C,D,E) = P(E|D)P(D|C)P(C|B)P(B|A)P(A)$$

$$P(E) = \sum_D \sum_C \sum_B \sum_A P(E|D)P(D|C)P(C|B)P(B|A)P(A)$$

$$= \sum_D P(E|D) \left[ \sum_C P(D|C) \left[ \sum_B P(C|B) \left[ \sum_A P(B|A)P(A) \right] \right] \right]$$

$$\begin{array}{c} \underbrace{\hspace{10em}}_{f(B)} \\ \underbrace{\hspace{12em}}_{g(C)} \\ \underbrace{\hspace{14em}}_{h(D)} \end{array}$$



# SP Algorithm as Free Energy Minimization

For a graphical model with potentials,  $\Psi_i(x_{C_i})$  on cliques  $x_{C_i}$ ,

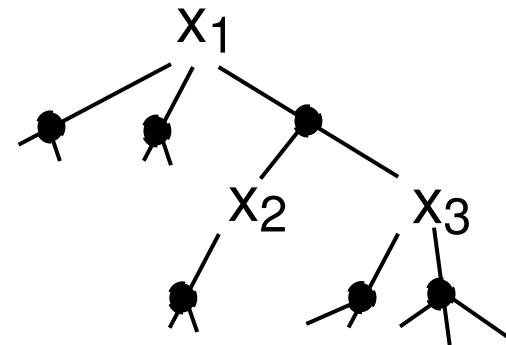
$$F = -\sum_i \sum_{x_{C_i}} Q(x_{C_i}) \log[\Psi_i(x_{C_i})] + \sum_x Q(x) \log[Q(x)]$$

- For a general Q-distribution, if  $Q(x_{C_i})$  are known, the 1<sup>st</sup> term is easy to compute
- Problem: 2<sup>nd</sup> term is generally intractable
- Idea: Represent  $Q(x)$  by  $Q(x_j)$ ,  $\forall j$ , and  $Q(x_{C_i})$ ,  $\forall i$ , and approximate the 2<sup>nd</sup> term

# Bethe approximation

- For a tree

$$Q(x) = \frac{\prod_i Q(x_{C_i})}{\prod_j Q(x_j)^{d_j-1}}$$



- $d_j = \text{degree of } x_j$
- Use this expression for the entropy, even when the graph is not a tree
- Note: For a tree, the minimizer of  $F$  gives  $Q(x) = P(x)$ : Exact inference

## Solving for Q

$$F = -\sum_i \sum_{x_{C_i}} Q(x_{C_i}) \log[\Psi_i(x_{C_i})] \\ + \sum_i \sum_{x_{C_i}} Q(x_{C_i}) \log[Q(x_{C_i})] - \sum_j (d_j - 1) \sum_{x_j} Q(x_j) \log[Q(x_j)]$$

Solve  $\partial D / \partial Q(x_{C_i}) = 0$  and  $\partial D / \partial Q(x_j) = 0$ ,  
subject to the constraints

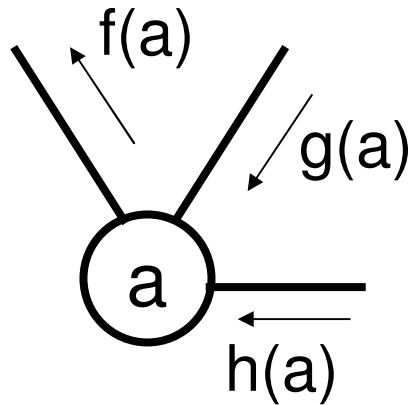
- $\sum_{x_j} Q(x_j) = 1$
- $\sum_{x_{C_i}} Q(x_{C_i}) = 1$
- For  $j \in C_i$ ,  $\sum_{x_{C_i} \setminus j} Q(x_{C_i}) = Q(x_j)$

Result (Yedidia et al 2001): Sum-product algorithm

# The sum-product algorithm

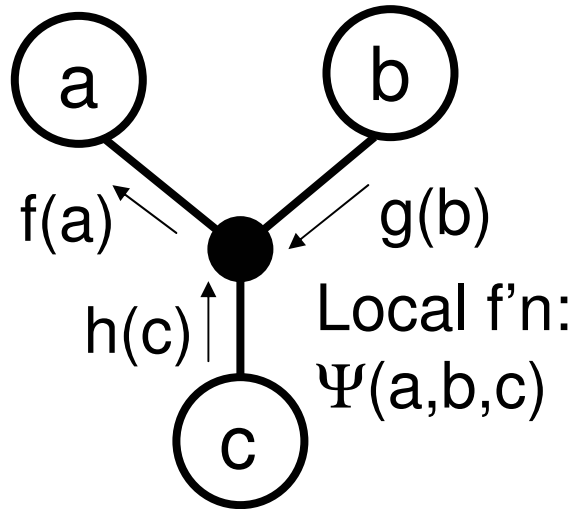
Each message is a function of its neighboring variable

Out of variable



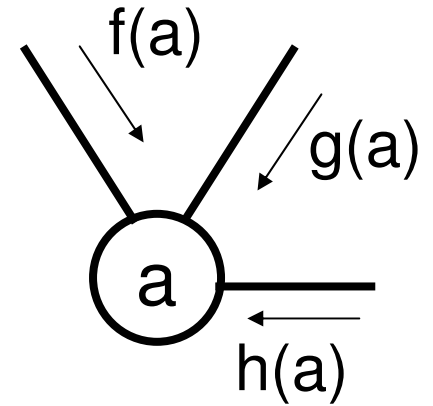
$$f(a) = g(a)h(a)$$

Out of function



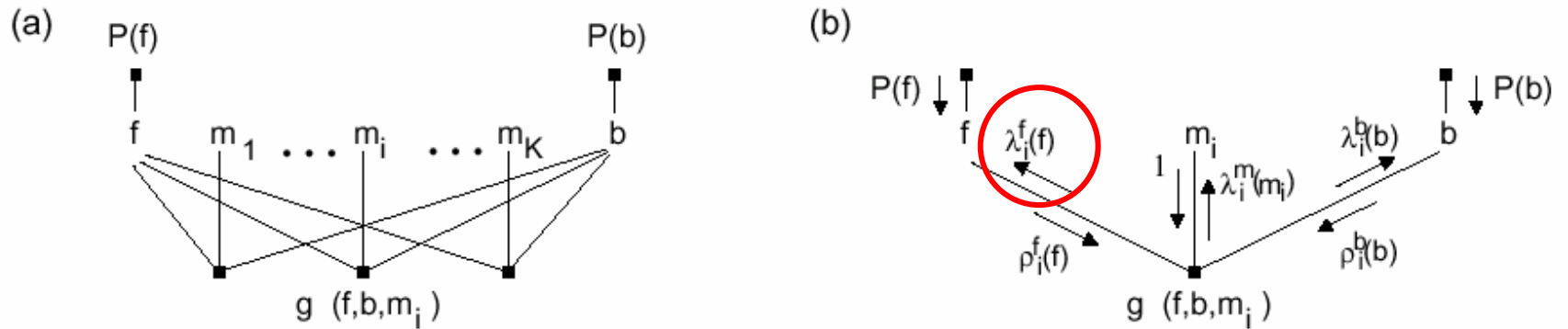
$$f(a) = \sum_b \sum_c \Psi(a,b,c)g(b)h(c)$$

Fusion



$$P(a) \approx f(a)g(a)h(a)$$

# Occlusion model: The sum-product algorithm



$$g_i(f, b, m_i) = P(z_i | m_i, f, b) P(m_i | f) =$$

$$N(z_i; \mu_{fi}, \psi_{fi})^{m_i} N(z_i; \mu_{bi}, \psi_{bi})^{1-m_i} \alpha_{fi}^{m_i} (1 - \alpha_{fi})^{1-m_i}.$$

$$\lambda_i^f(f) \leftarrow \alpha_{fi} N(z_i; \mu_{fi}, \psi_{fi}) + (1 - \alpha_{fi}) \sum_b N(z_i; \mu_{bi}, \psi_{bi}) \rho_i^b(b).$$

$$\lambda_i^f(f) \leftarrow \lambda_i^f(f) / (\sum_f \lambda_i^f(f))$$

# Occlusion model: The sum-product algorithm



$$g_i(f, b, m_i) = P(z_i | m_i, f, b) P(m_i | f) =$$

$$N(z_i; \mu_{fi}, \psi_{fi})^{m_i} N(z_i; \mu_{bi}, \psi_{bi})^{1-m_i} \alpha_{fi}^{m_i} (1 - \alpha_{fi})^{1-m_i}.$$

$$\rho_i^f(f) \leftarrow P(f) \prod_{j \neq i} \lambda_j^f(f),$$

$$\rho_i^f(f) \leftarrow \rho_i^f(f) / (\sum_f \rho_i^f(f))$$

# Occlusion model: The sum-product algorithm



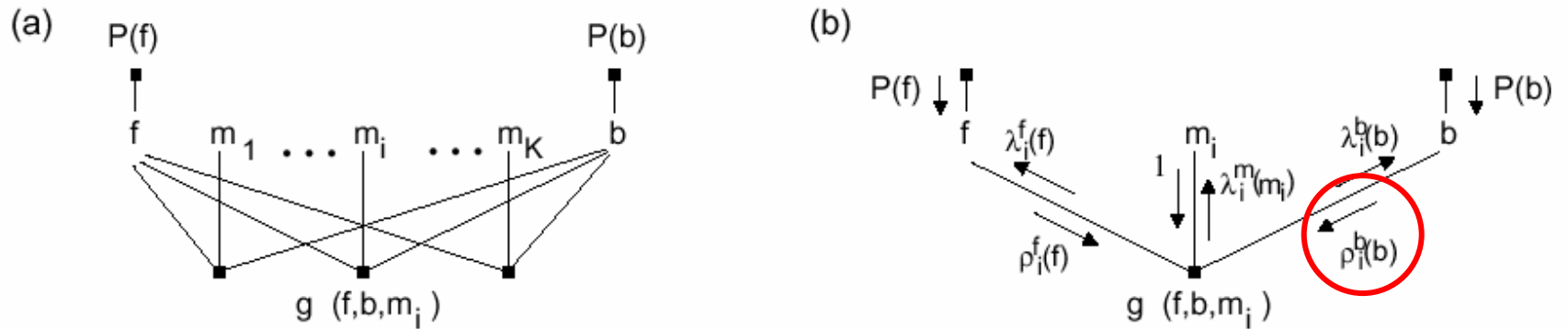
$$g_i(f, b, m_i) = P(z_i | m_i, f, b) P(m_i | f) =$$

$$N(z_i; \mu_{fi}, \psi_{fi})^{m_i} N(z_i; \mu_{bi}, \psi_{bi})^{1-m_i} \alpha_{fi}^{m_i} (1 - \alpha_{fi})^{1-m_i}.$$

$$\lambda_i^b(b) \leftarrow \sum_f \sum_{m_i} g_i(f, b, m_i) \cdot 1 \cdot \rho_i^f(f)$$

$$\lambda_i^b(b) \leftarrow \lambda_i^b(b) / \left( \sum_b \lambda_i^b(b) \right)$$

# Occlusion model: The sum-product algorithm



$$g_i(f, b, m_i) = P(z_i | m_i, f, b) P(m_i | f) =$$

$$N(z_i; \mu_{fi}, \psi_{fi})^{m_i} N(z_i; \mu_{bi}, \psi_{bi})^{1-m_i} \alpha_{fi}^{m_i} (1 - \alpha_{fi})^{1-m_i}.$$

$$\rho_i^b(b) \leftarrow P(b) \prod_{j \neq i} \lambda_j^b(b),$$

$$\rho_i^b(b) \leftarrow \rho_i^b(b) / (\sum_b \rho_i^b(b))$$



# Occlusion model: The sum-product algorithm



$$g_i(f, b, m_i) = P(z_i | m_i, f, b) P(m_i | f) =$$

$$N(z_i; \mu_{fi}, \psi_{fi})^{m_i} N(z_i; \mu_{bi}, \psi_{bi})^{1-m_i} \alpha_{fi}^{m_i} (1 - \alpha_{fi})^{1-m_i}.$$

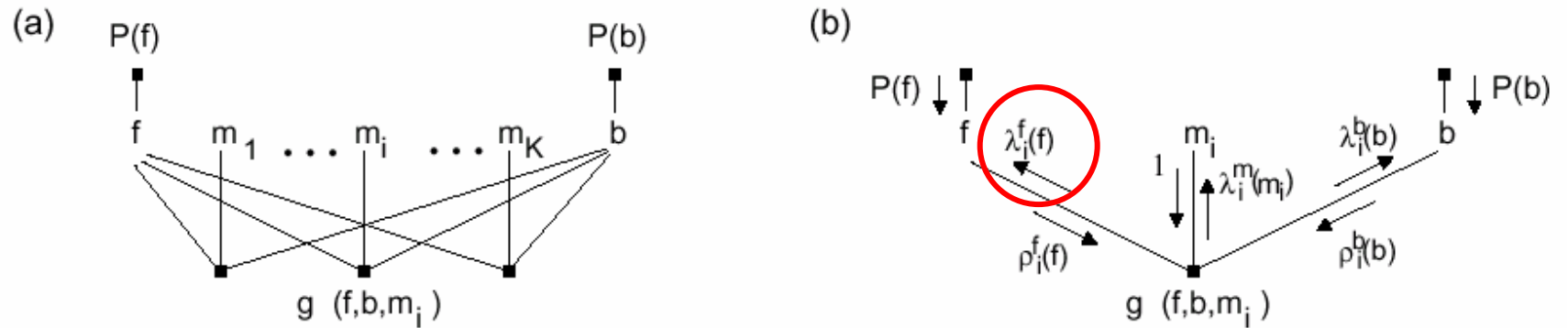
$$\lambda_i^m(m_i) \leftarrow \sum_f \sum_b g_i(f, b, m_i) \cdot \rho_i^f(f) \cdot \rho_i^b(b) :$$

$$\lambda_i^m(1) \leftarrow \sum N(z_i; \mu_{fi}, \psi_{fi}) \alpha_{fi} \rho_i^f(f),$$

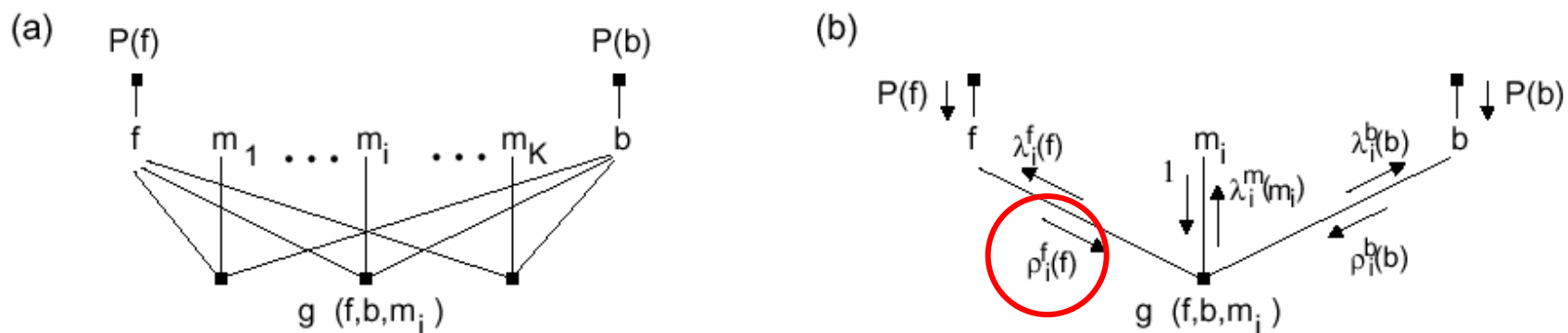
$$\lambda_i^m(0) \leftarrow \left( \sum_b N(z_i; \mu_{bi}, \psi_{bi}) \rho_i^b(b) \right) \left( \sum_f (1 - \alpha_{fi}) \rho_i^f(f) \right).$$

$$\lambda_i^m(m_i) \leftarrow \lambda_i^m(m_i) / (\lambda_i^m(0) + \lambda_i^m(1))$$

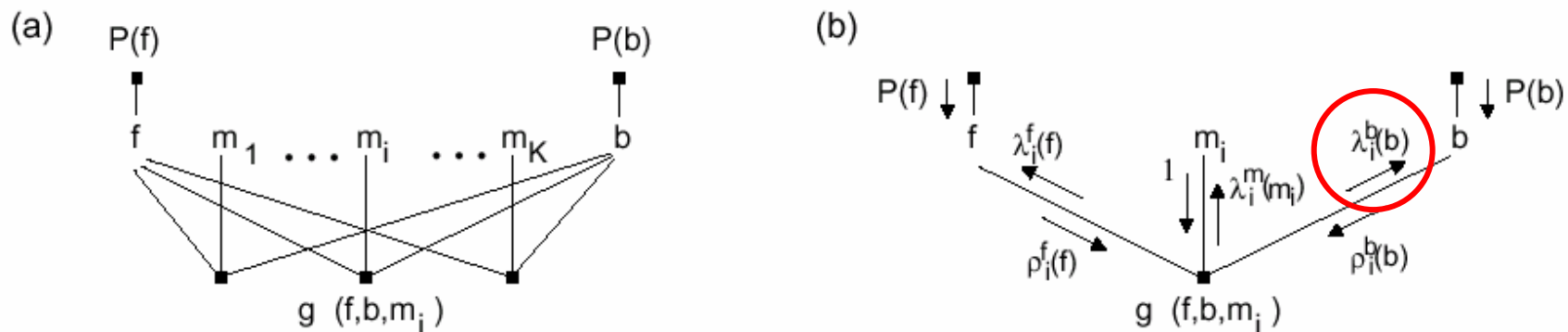
# Occlusion model: The sum-product algorithm



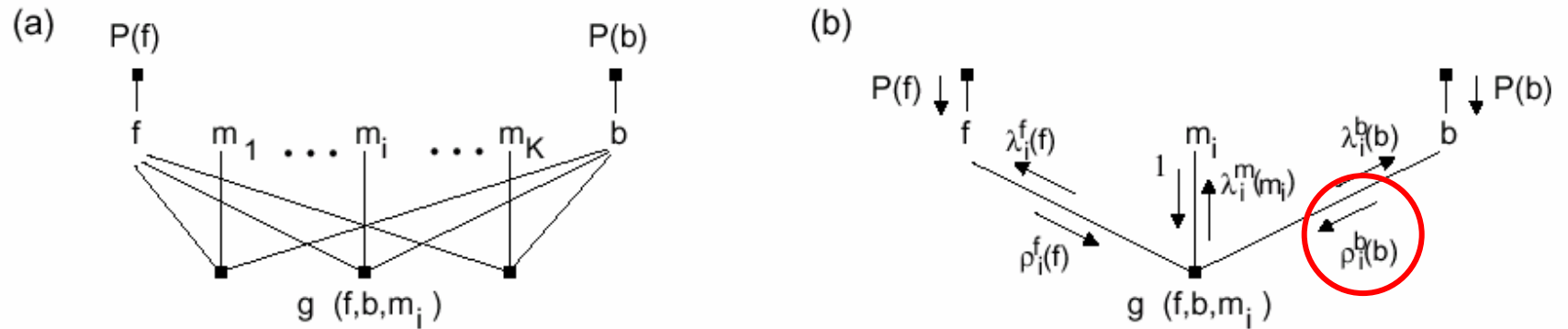
# Occlusion model: The sum-product algorithm



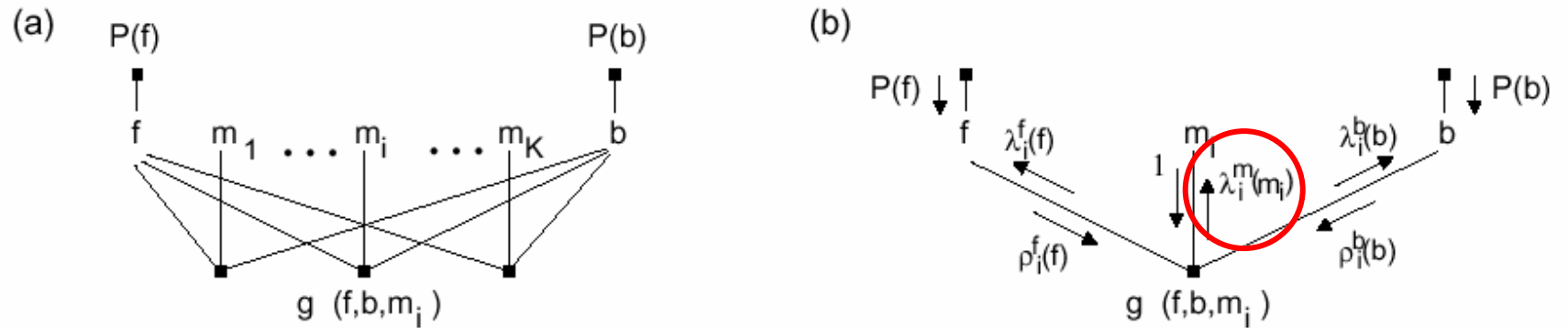
# Occlusion model: The sum-product algorithm



# Occlusion model: The sum-product algorithm



# Occlusion model: The sum-product algorithm



# Gibbs sampling

- The posterior  $P(h_1, h_2, \dots, h_K | v)$  is not tractable, but often, the conditionals,  $P(h_i | h \setminus h_i, v)$ , can be sampled from
- Gibbs sampling:

**Initialization.** Pick values for all hidden variables (randomly, or cleverly).

**Sampling Step.** Choose a variable  $h_i$  at random or in order, and then sample it from  $P(h_i | h \setminus h_i)$ .

**Repeat for a fixed number of iterations or until convergence.**

# Gibbs sampling as free energy minimization

- Imagine running an ensemble of Gibbs chains in parallel
- Let  $Q^n(h)$  describe the distribution of the  $h$ 's at step  $n$  of the sampling procedure

- Suppose at step  $n$ , we sample  $h_i$  in every chain:

$$Q^{n+1}(h) = Q^n(h \setminus h_i)P(h_i|h \setminus h_i, v)$$

- Substituting these into the expression for  $F$ , we find that  $F^{n+1} \leq F^n$



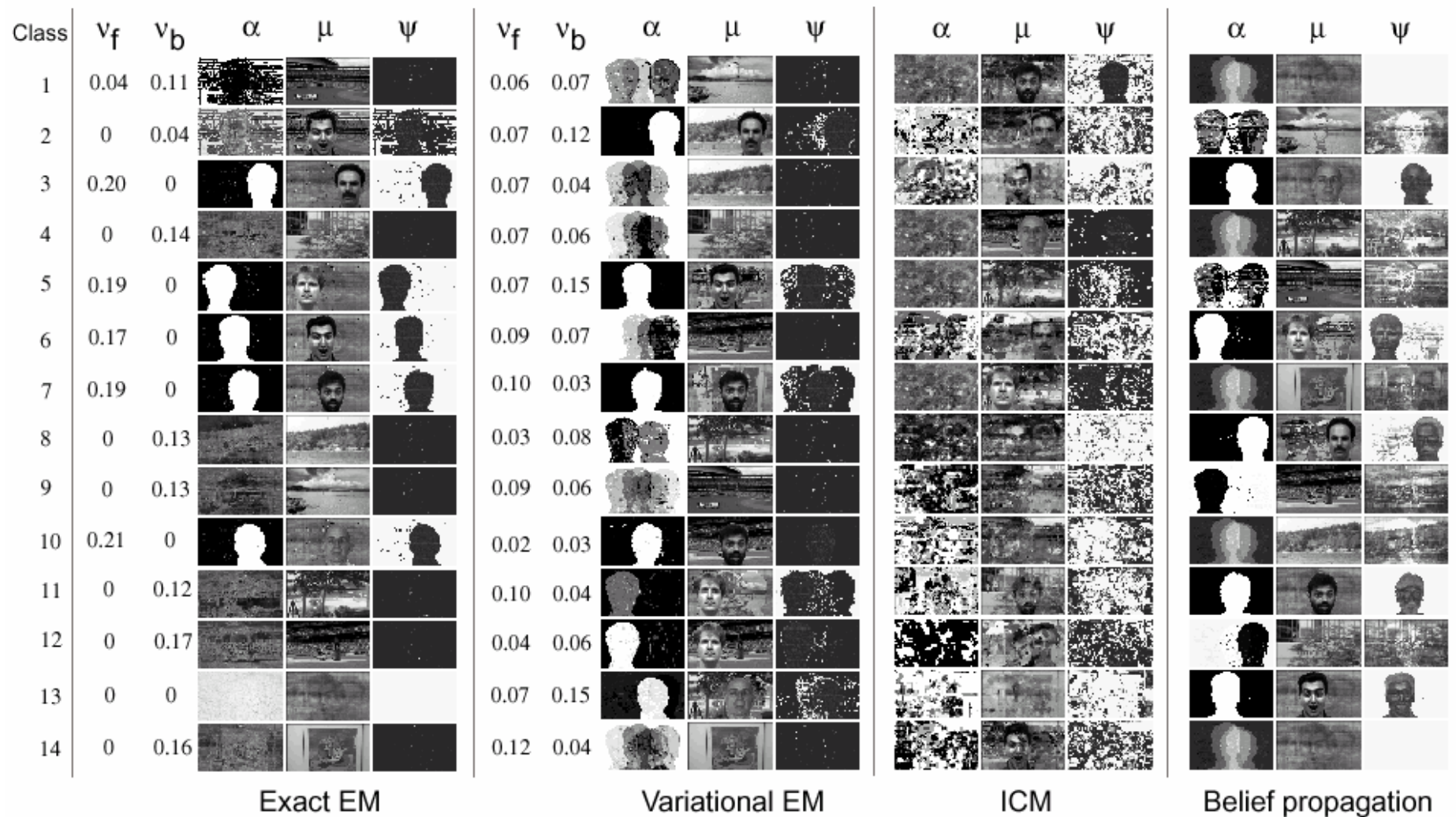
# Discussion of Algorithms

B. J. FREY & N. JOJIC

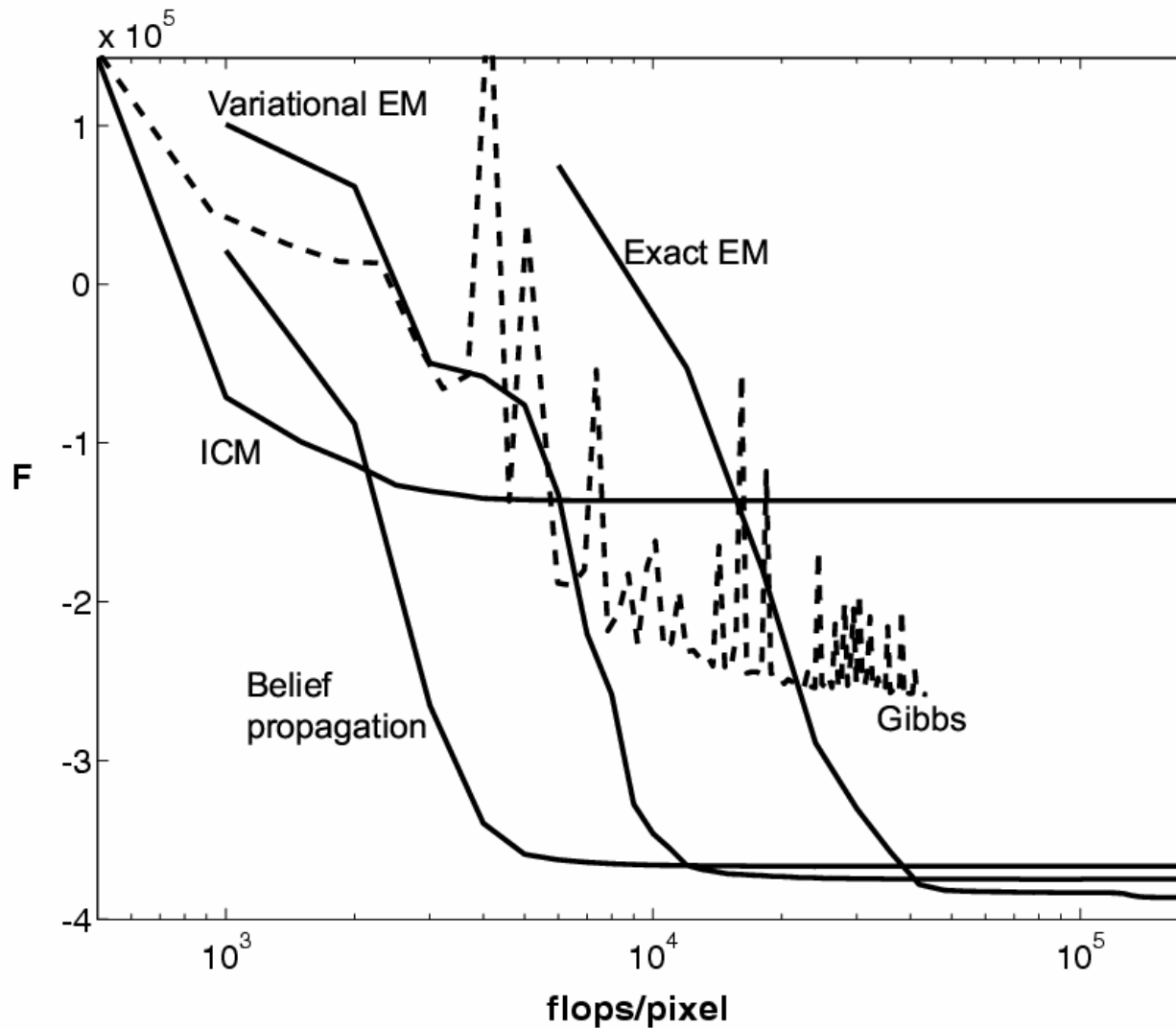
# Comparison of algorithms for occlusion model

Method	Update for mask variables	Complexity
Exact inference (used in EM)	$\frac{Q(m_i=1 b,f)}{Q(m_i=0 b,f)} \leftarrow \frac{\alpha_{f_i} \mathcal{N}(z_i; \mu_{f_i}, \psi_{f_i})}{(1-\alpha_{f_i}) \mathcal{N}(z_i; \mu_{b_i}, \psi_{b_i})}$	$J^2K$
Iterative conditional modes	$m_i \leftarrow \begin{cases} 1, & \text{if } \frac{\alpha_{f_i} \mathcal{N}(z_i; \mu_{f_i}, \psi_{f_i})}{(1-\alpha_{f_i}) \mathcal{N}(z_i; \mu_{b_i}, \psi_{b_i})} > 1 \\ 0, & \text{otherwise} \end{cases}$	$K$
Gibbs sampling	$m_i \leftarrow \text{sample}_{m_i} \left\{ \begin{array}{ll} \alpha_{f_i} \mathcal{N}(z_i; \mu_{f_i}, \psi_{f_i}) & \text{if } m_i = 1 \\ (1-\alpha_{f_i}) \mathcal{N}(z_i; \mu_{b_i}, \psi_{b_i}) & \text{if } m_i = 0 \end{array} \right\}$	$K$
Fully-factorized variational	$\frac{Q(m_i=1)}{Q(m_i=0)} \leftarrow \frac{\prod_f (\alpha_{f_i} \mathcal{N}(z_i; \mu_{f_i}, \psi_{f_i}))^{Q(f)}}{(\prod_f (1-\alpha_{f_i})^{Q(f)}) (\prod_b \mathcal{N}(z_i; \mu_{b_i}, \psi_{b_i})^{Q(b)})}$	$JK$
Structured variational	$\frac{Q(m_i=1 f)}{Q(m_i=0 f)} \leftarrow \frac{\alpha_{f_i} \mathcal{N}(z_i; \mu_{f_i}, \psi_{f_i})}{(1-\alpha_{f_i}) \prod_b \mathcal{N}(z_i; \mu_{b_i}, \psi_{b_i})^{Q(b)}}$	$JK$
Sum-product algorithm	$\frac{Q(m_i=1)}{Q(m_i=0)} \leftarrow \frac{\sum_f \rho_i^f(f) \alpha_{f_i} \mathcal{N}(z_i; \mu_{f_i}, \psi_{f_i})}{(\sum_f \rho_i^f(f) (1-\alpha_{f_i})) (\sum_b \rho_i^b(b) \mathcal{N}(z_i; \mu_{b_i}, \psi_{b_i}))}$	$JK$

# Experimental results for occlusion model



# Speed of convergence



# Recall: From ICM to Gibbs sampling

## Sample

- For  $t = 1, \dots, T$

$$\{ f^{(t)} \leftarrow \operatorname{argmax}_{f^{(t)}} [\pi_{f^{(t)}} \prod_{i: m_i^{(t)}=1} \mathcal{N}(z_i^{(t)}; \mu_{f^{(t)}i}, \psi_{f^{(t)}i})]$$

$$\{ b^{(t)} \leftarrow \operatorname{argmax}_{b^{(t)}} [\pi_{b^{(t)}} \prod_{i: m_i^{(t)}=0} \mathcal{N}(z_i^{(t)}; \mu_{b^{(t)}i}, \psi_{b^{(t)}i})]$$

$$\{ \text{For } i = 1, \dots, K: m_i^{(t)} \leftarrow \begin{cases} 1 & \text{if } \alpha_{f^{(t)}i} \mathcal{N}(z_i^{(t)}; \mu_{f^{(t)}i}, \psi_{f^{(t)}i}) > (1 - \alpha_{f^{(t)}i}) \mathcal{N}(z_i^{(t)}; \mu_{b^{(t)}i}, \psi_{b^{(t)}i}) \\ 0 & \text{otherwise} \end{cases}$$

- For  $j = 1, \dots, J$

$$\{ \pi_j \leftarrow (\sum_{t=1}^T [f^{(t)} = j] + \sum_{t=1}^T [b^{(t)} = j]) / 2T$$

- For  $j = 1, \dots, J$ , for  $i = 1, \dots, K$

$$\{ \alpha_{ji} \leftarrow (\sum_{t=1}^T [f^{(t)} = j] m_i^{(t)}) / (\sum_{t=1}^T [f^{(t)} = j])$$

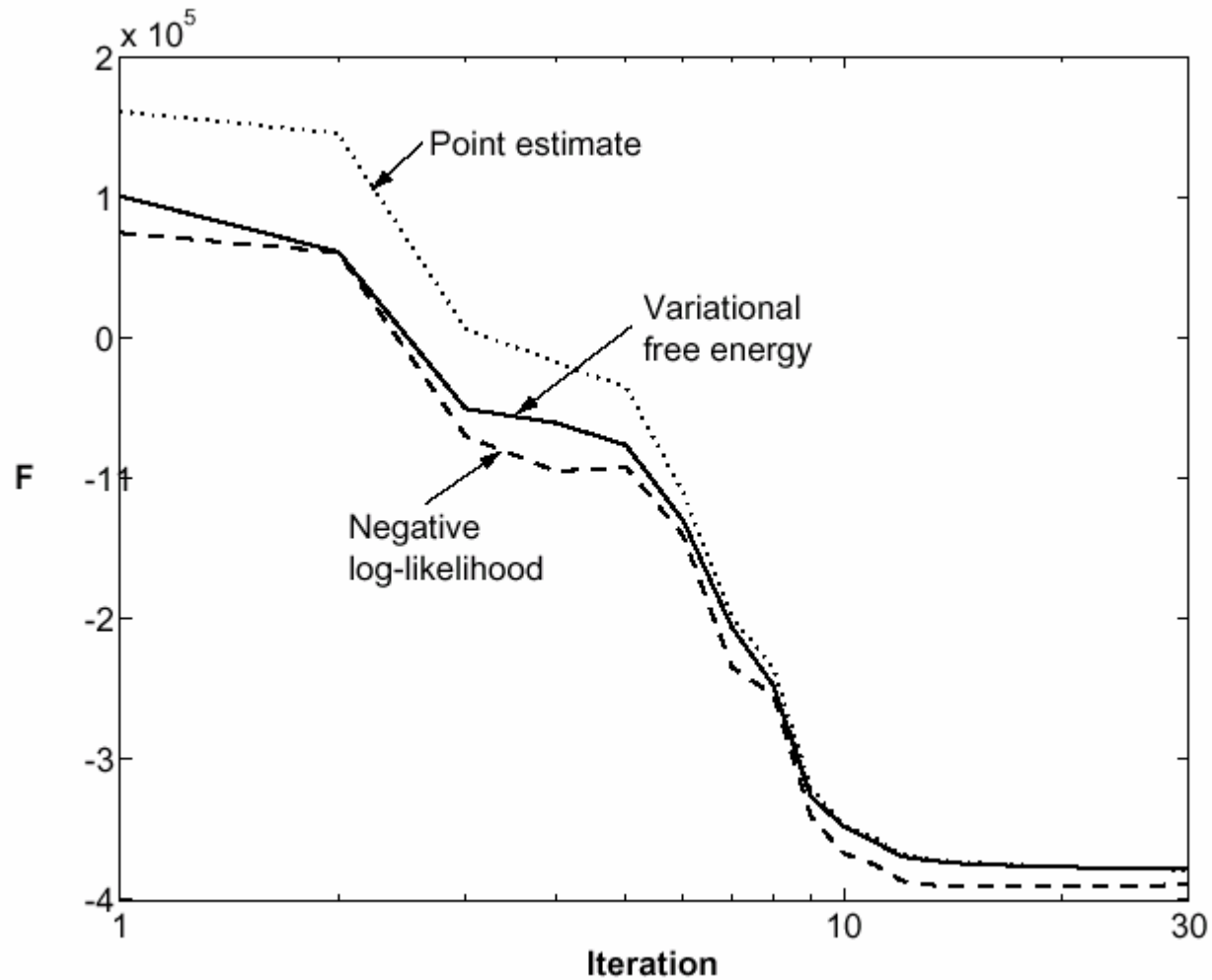
$$\{ \mu_{ji} \leftarrow (\sum_{t=1}^T [f^{(t)} = j \text{ or } b^{(t)} = j] z_i^{(t)}) / (\sum_{t=1}^T [f^{(t)} = j \text{ or } b^{(t)} = j])$$

$$\{ \psi_{ji} \leftarrow (\sum_{t=1}^T [f^{(t)} = j \text{ or } b^{(t)} = j] (z_i^{(t)} - \mu_{ji})^2) / (\sum_{t=1}^T [f^{(t)} = j \text{ or } b^{(t)} = j])$$

Here, the Iverson notation is used where [True] = 1 and [False] = 0.

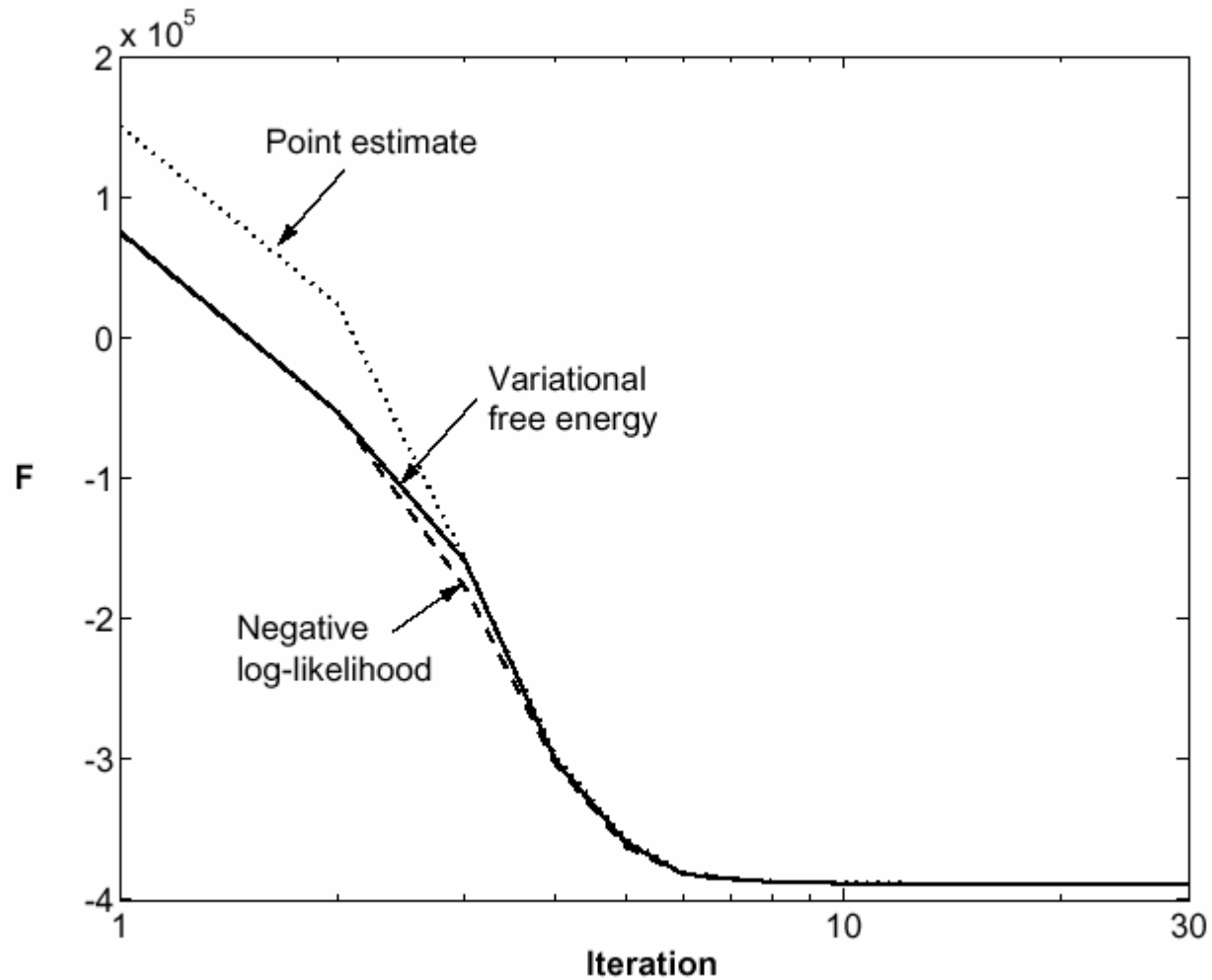
# How tight are the bounds?

Various bounds during variational learning



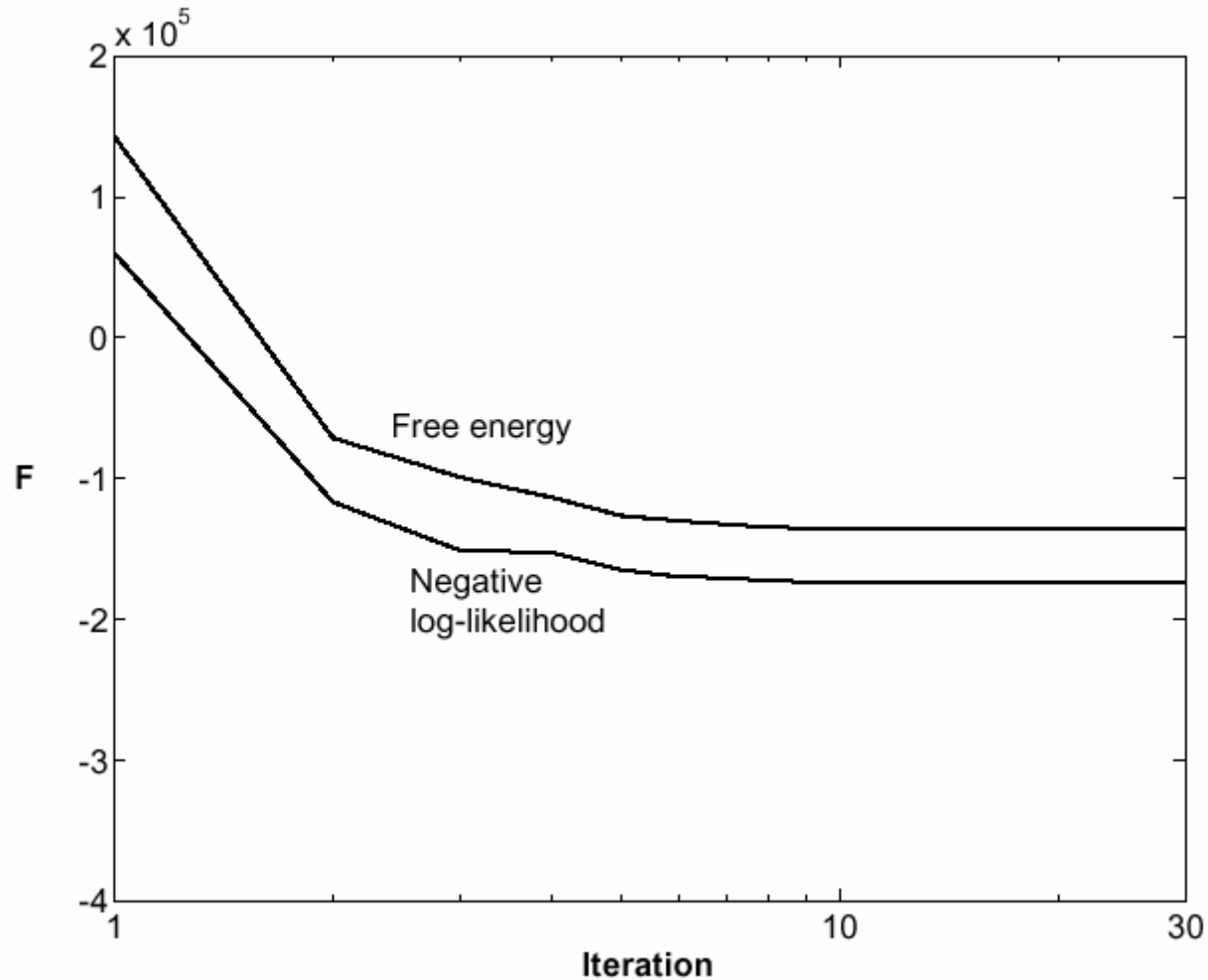
# How tight are the bounds?

Various bounds during EM learning



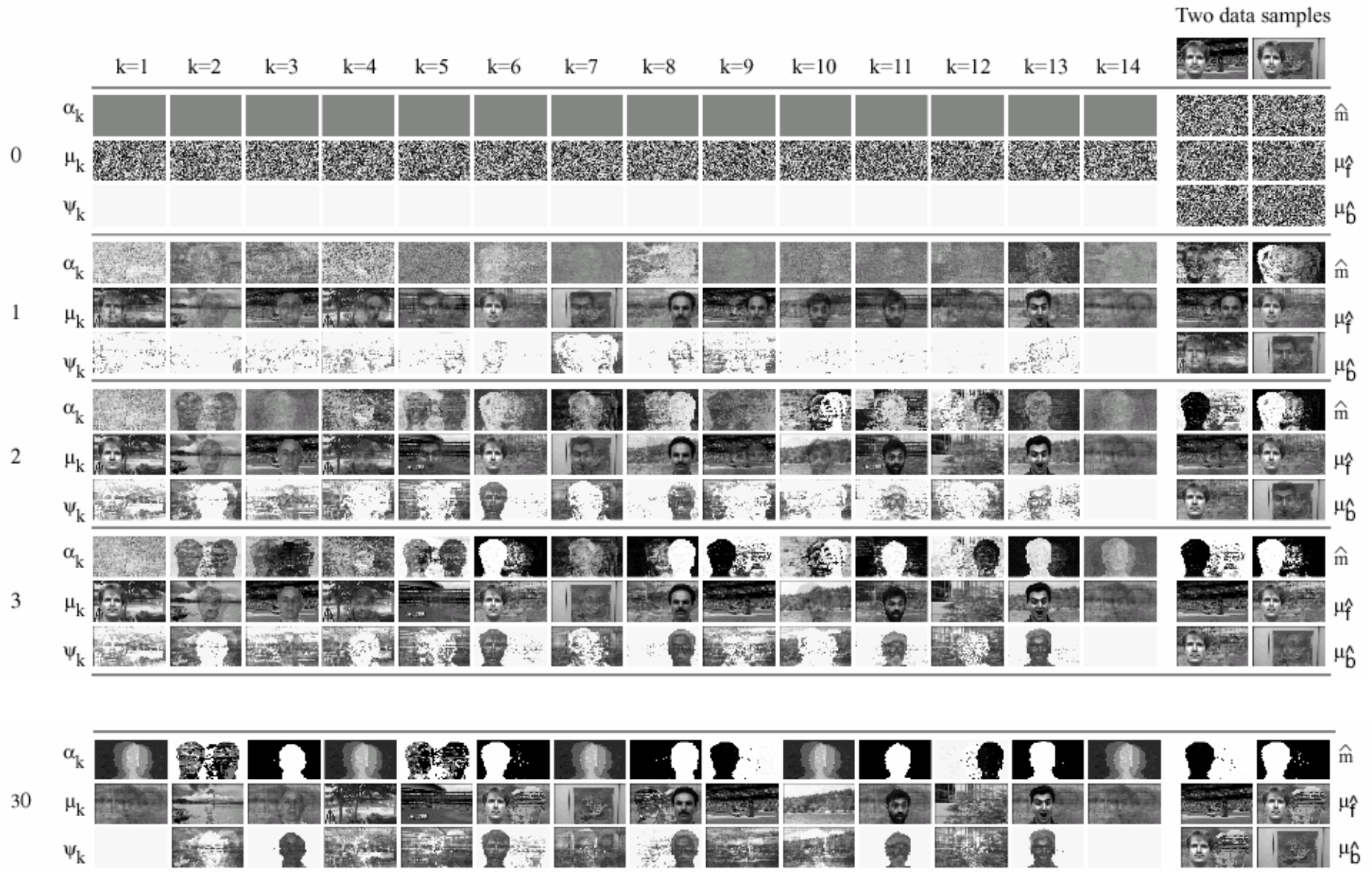
# How tight are the bounds?

Various bounds during ICM learning





# Iterative model optimization (sum-product algorithm)



# Improving approximate methods with additional EM iterations



# Wrap-up

# The Art

- Representation is critical
- Structure versus forms of probability functions
- What computation should the graph be used for (eg, recursive updates or coordinated optimization)?
- Which is better: Exact inference in inaccurate models or approximate inference in accurate models?
- ...

# Application: Cardboard cut-out scene analysis

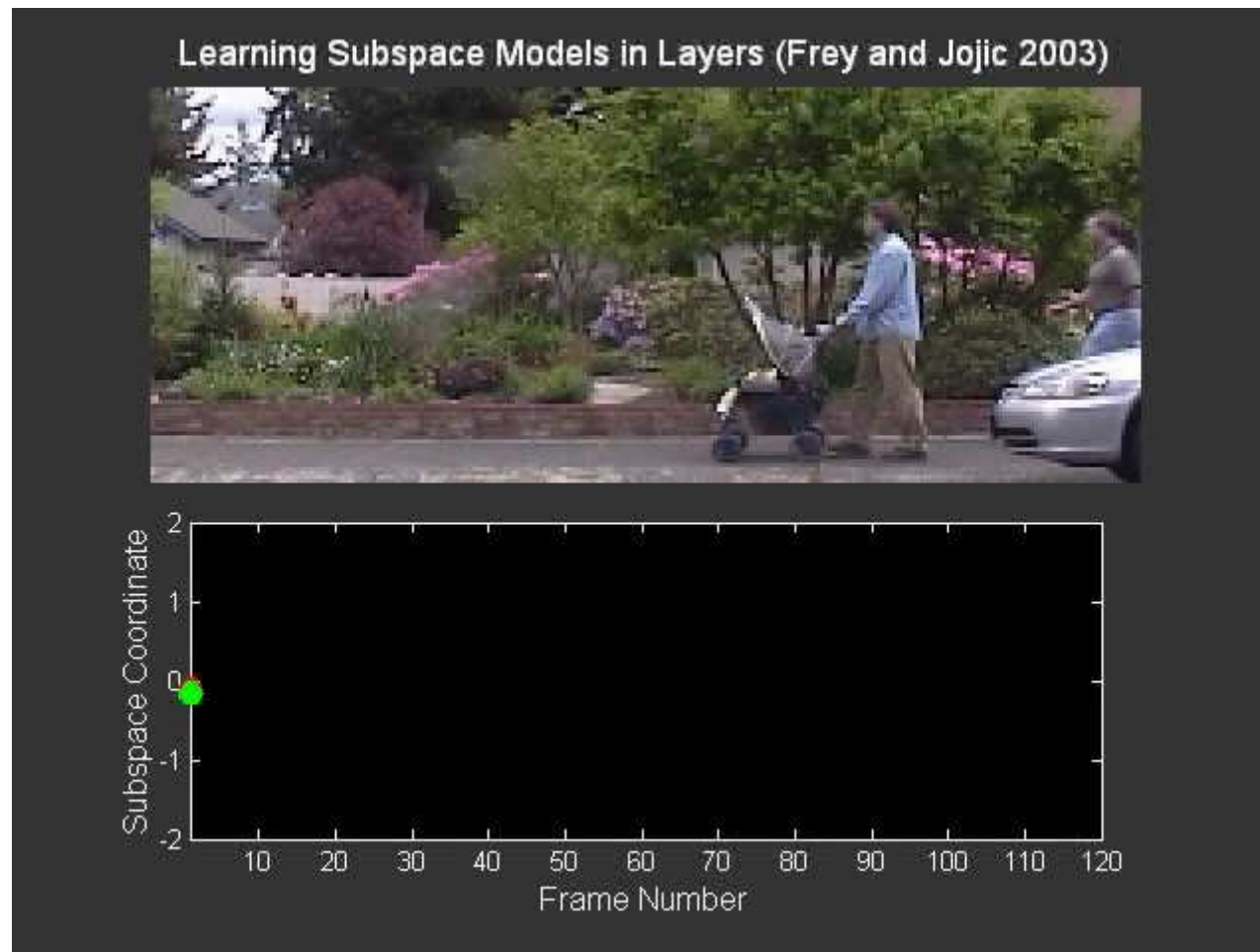
Jojic and Frey, CVPR 2001



B. J. FREY & N. JOJIC

# Application: Subspace models of occluding objects in 3-D scenes

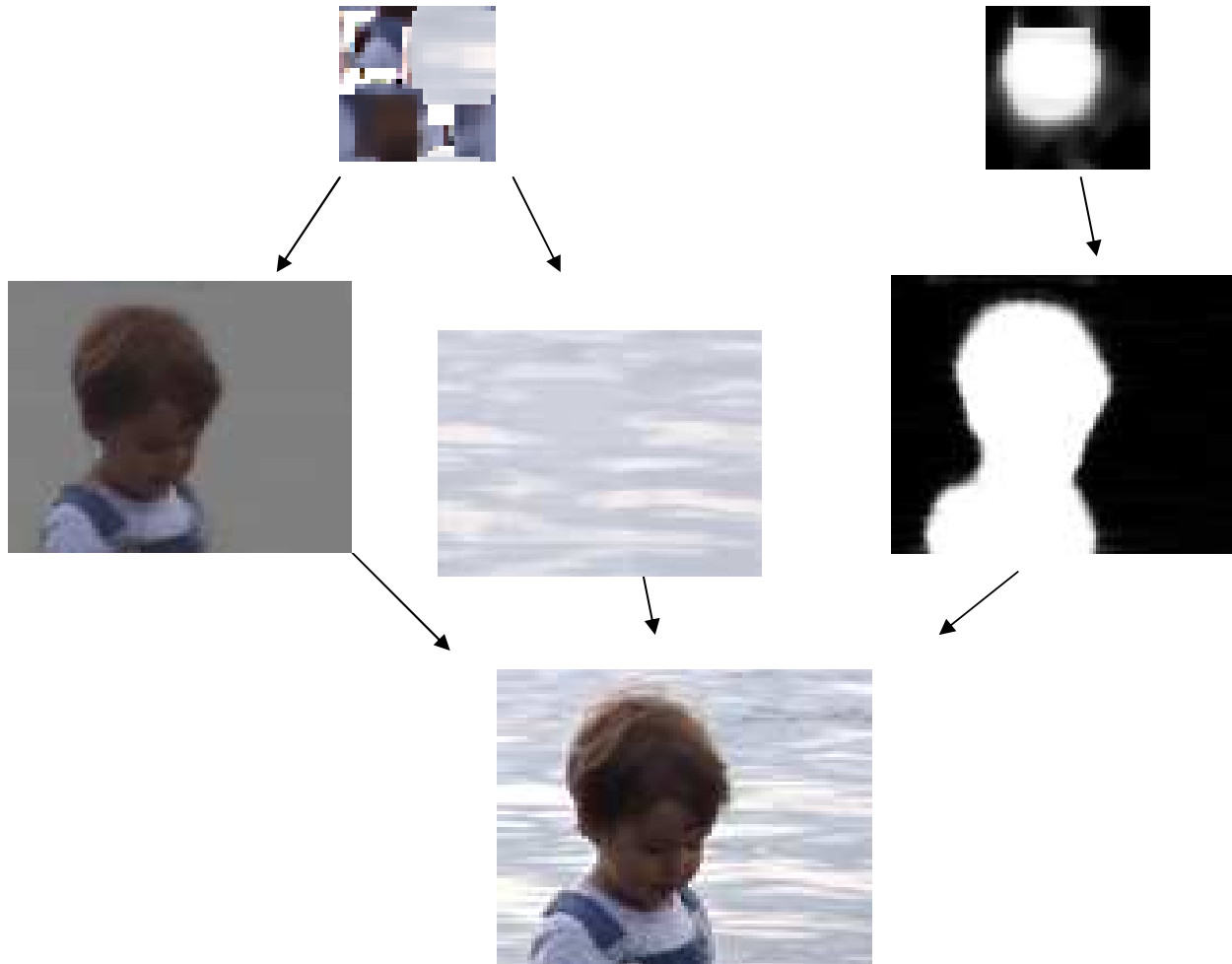
Frey, Jojic and Kannan 2003



B. J. FREY & N. JOJIC

# Application: Scene interpretation from single images

Jojic, Frey and Kannan 2003



- Those are the basics
- Go forth, model, perform approximate inference, and have fun!