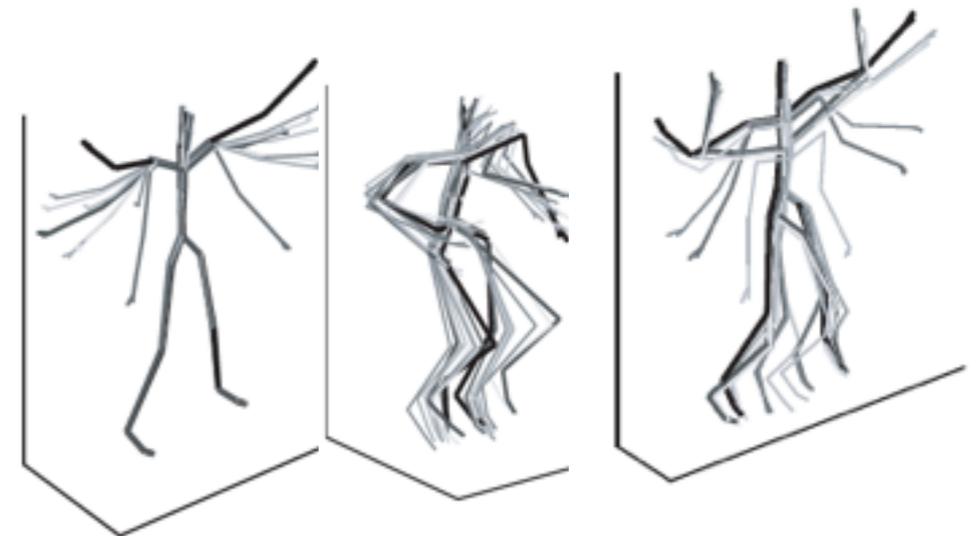
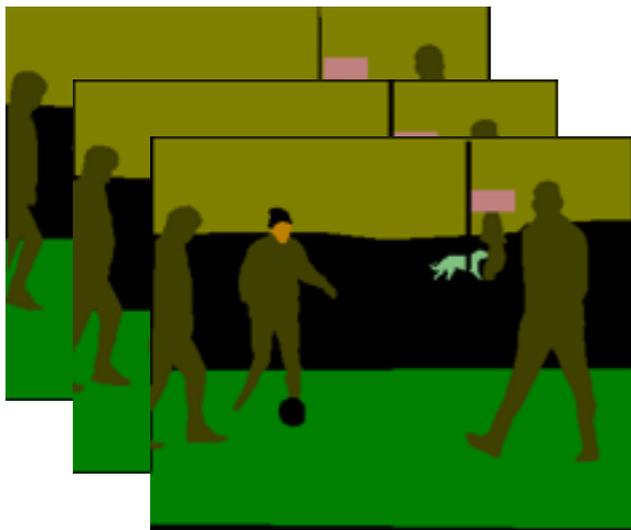
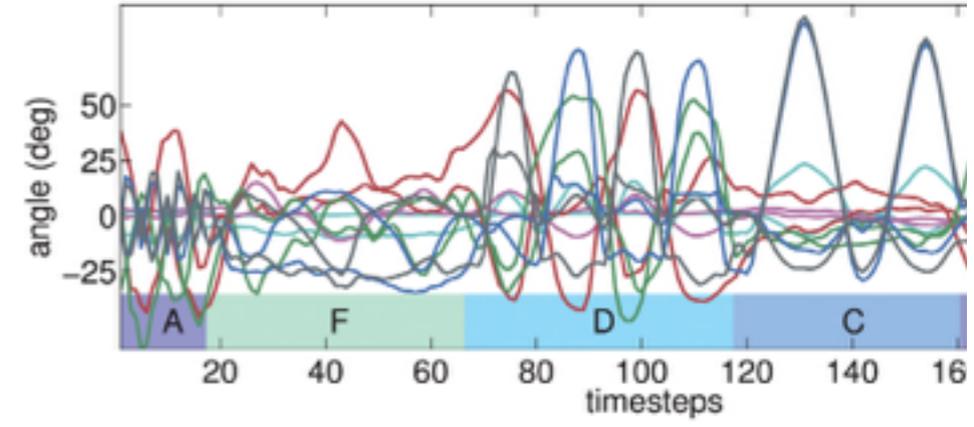


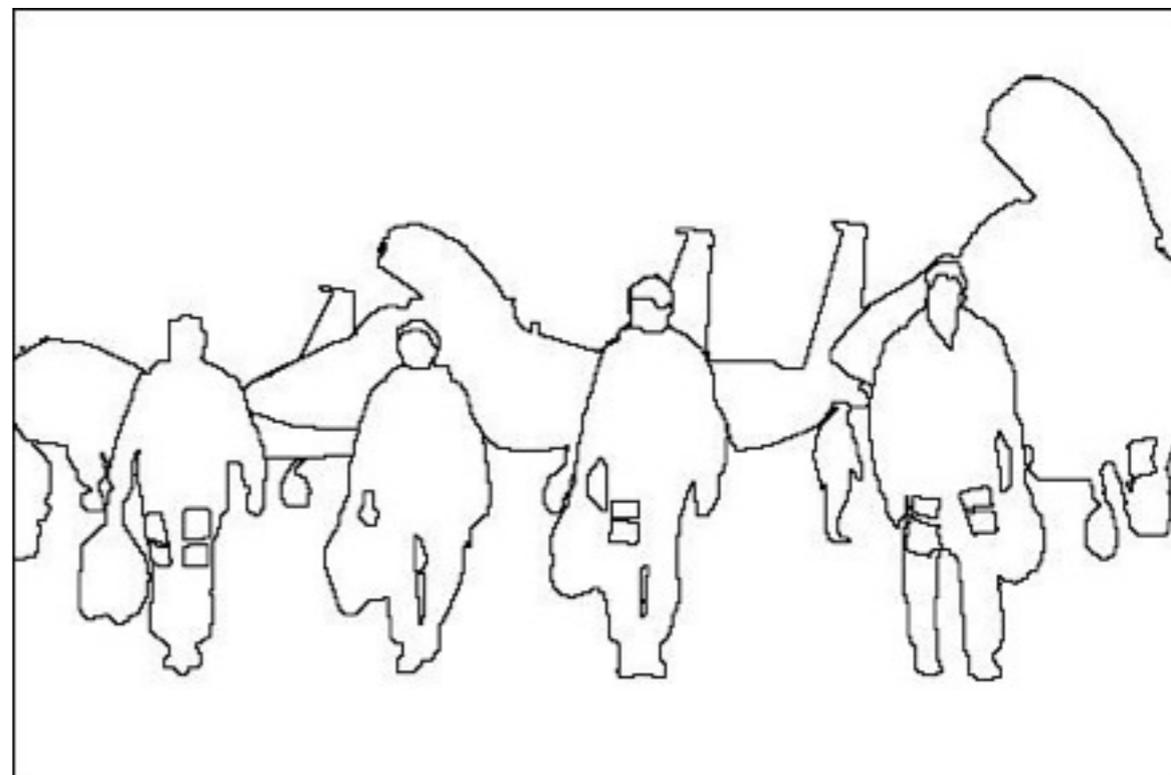
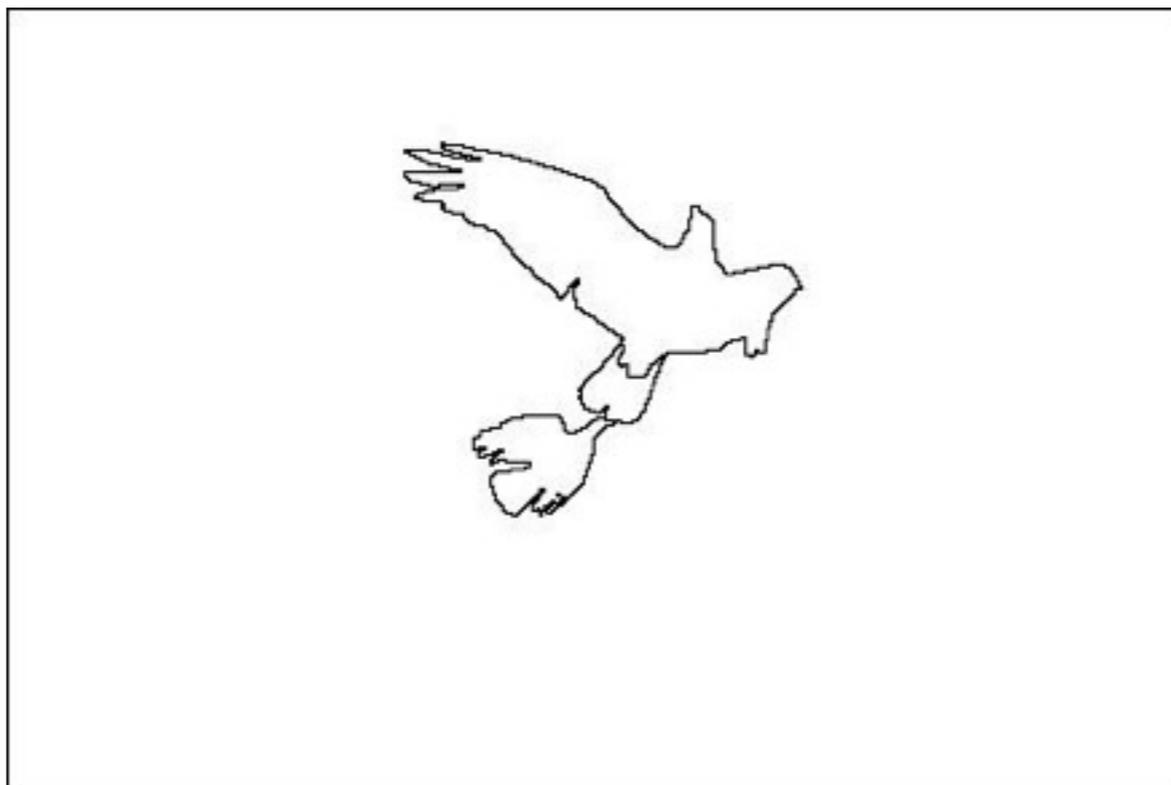
Bayesian Nonparametric Discovery of Layers and Parts from Scenes and Objects

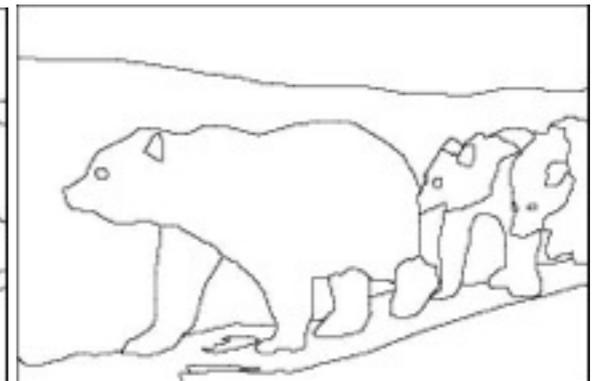
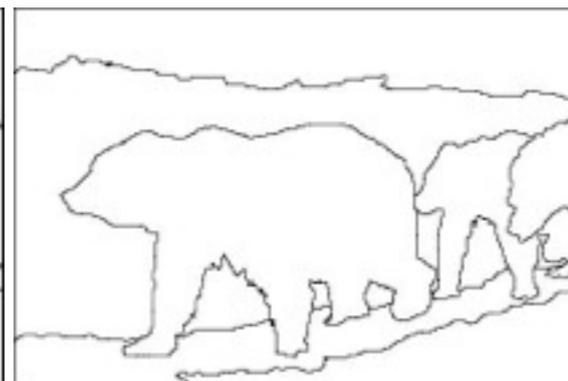
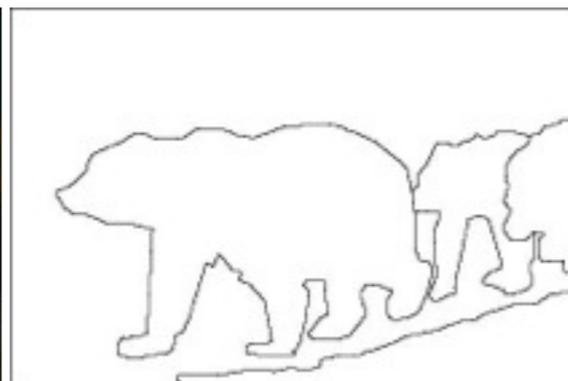
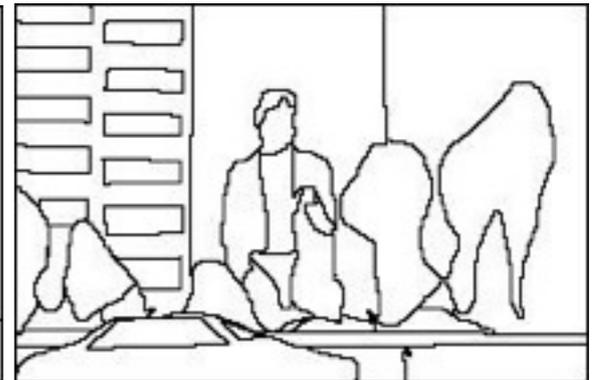
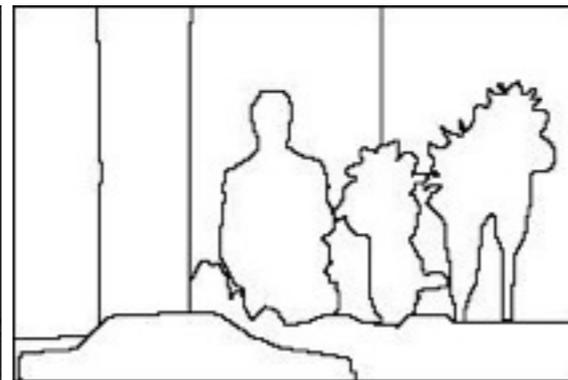
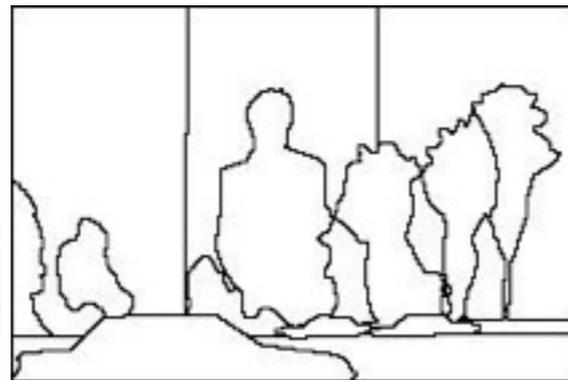
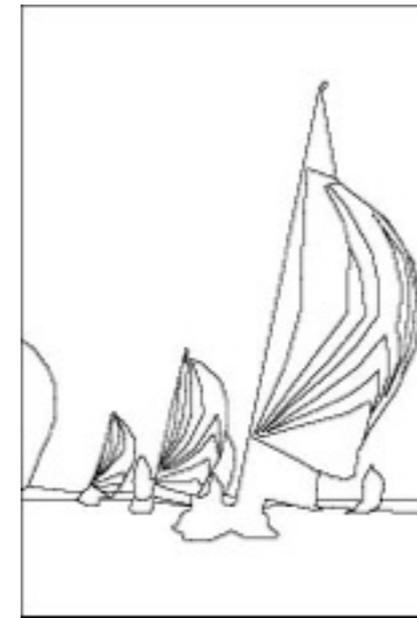
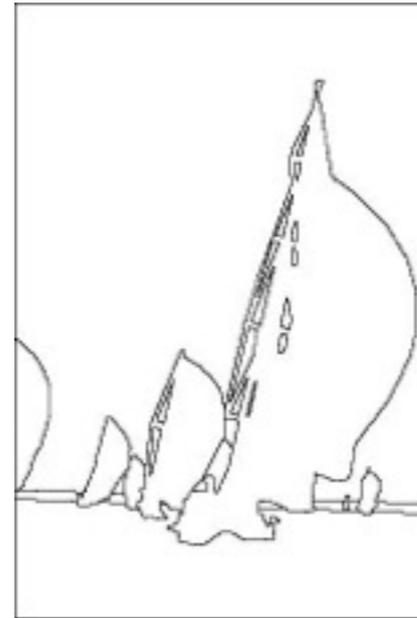
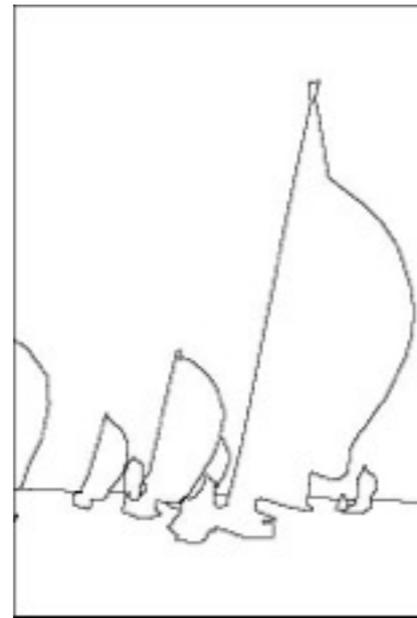
Soumya Ghosh

Advisor: Erik Sudderth

Committee: Michael Black and James Hays





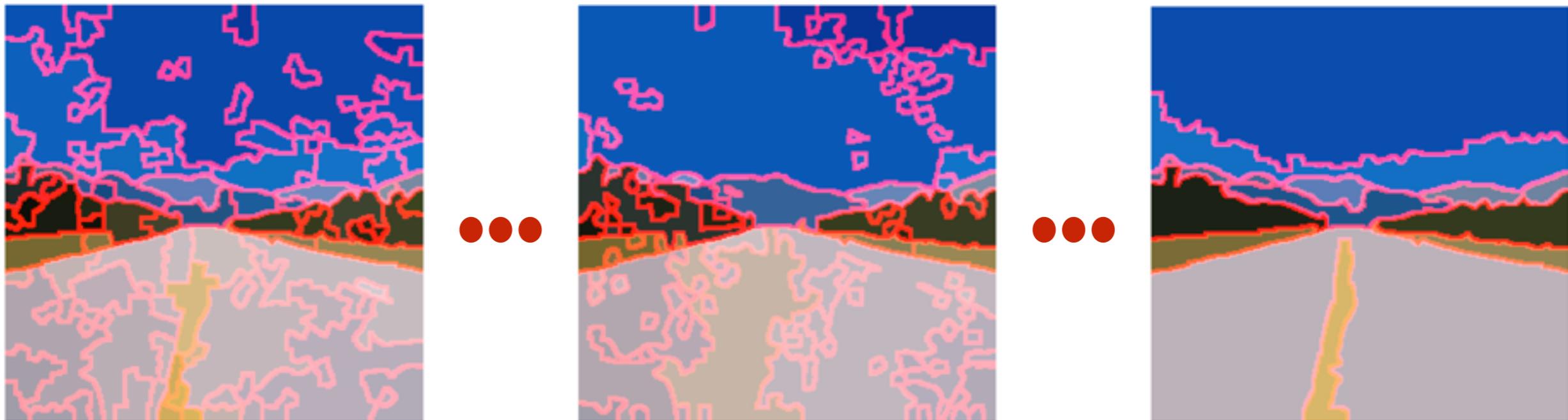


Human Segmentations

Model Desiderata

- Automatic **model selection** - adapt to variability in image/video/object complexity
- Manage **uncertainty** - retain a distribution over possible explanations
- Model **spatial** and **temporal correlations**
- **Learn** from **human** explanations

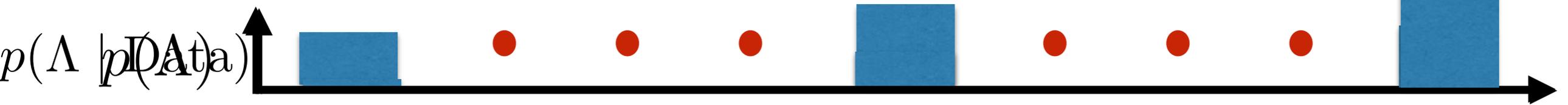
Adapting to complexity: Distributions over partitions



Λ_1

Λ_2

Λ_3

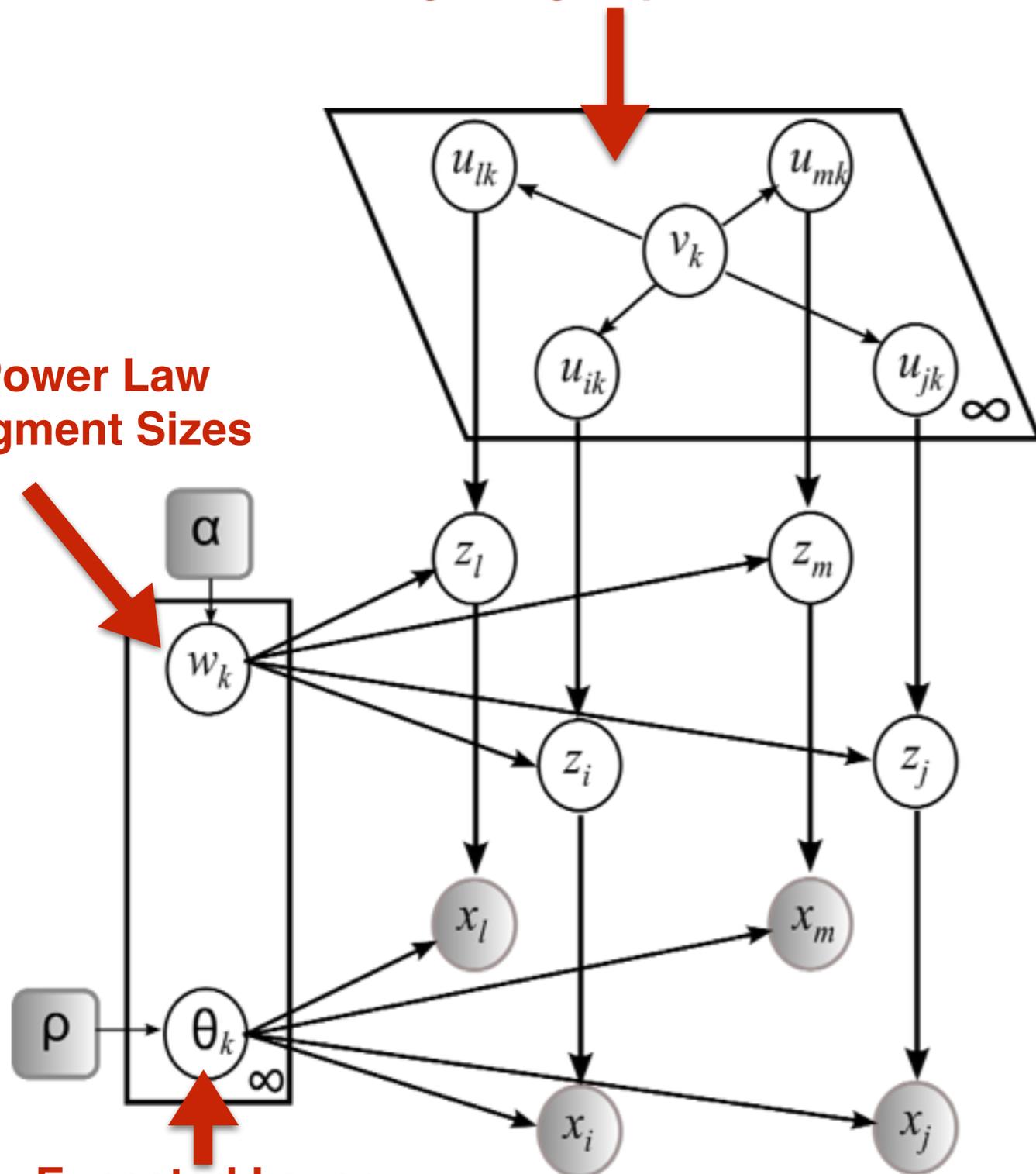


$$\Lambda^* \sim p(\Lambda | \text{Data})$$

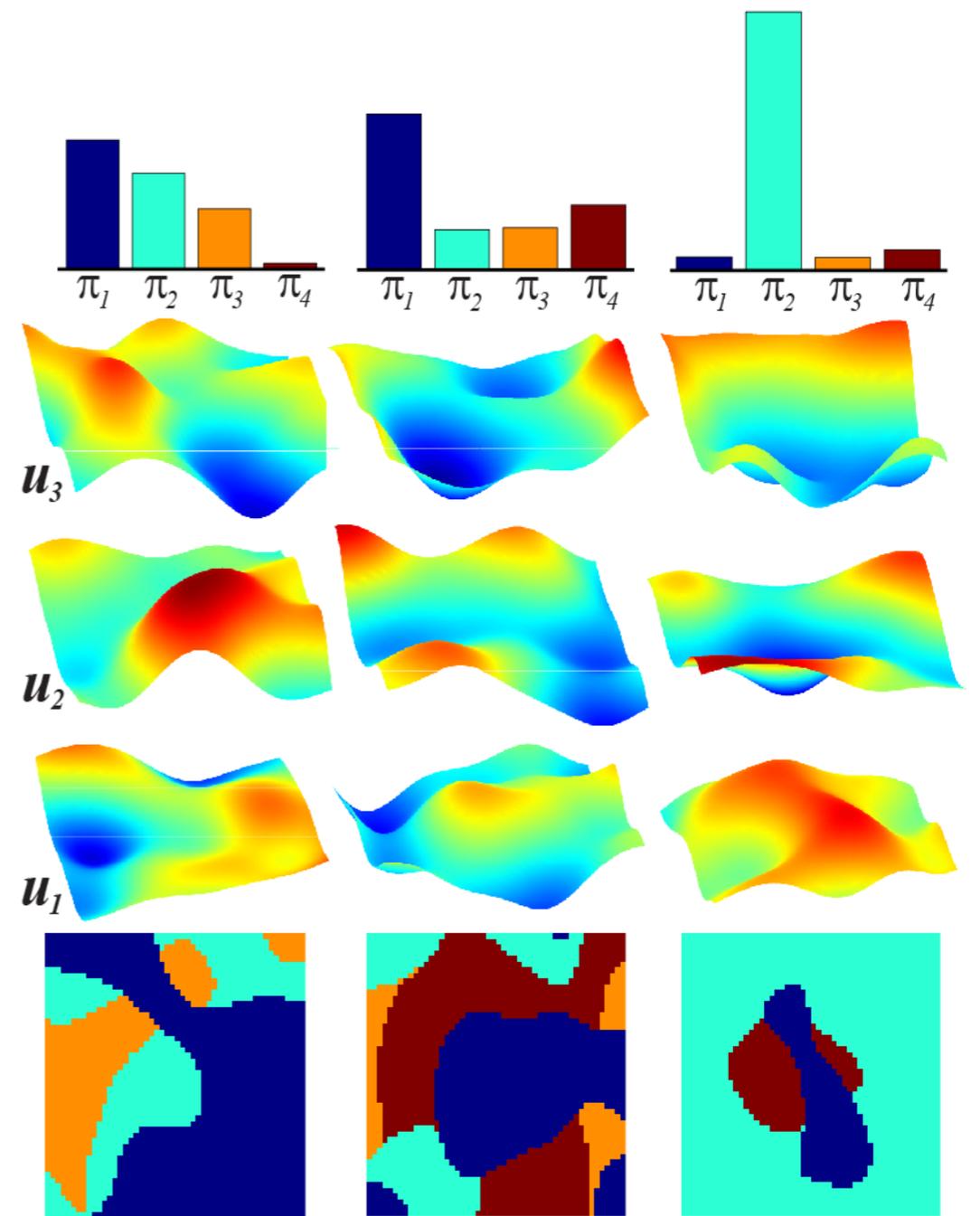
Spatially Coupled PY Processes

Model Long Range Spatial Correlations

Power Law Segment Sizes

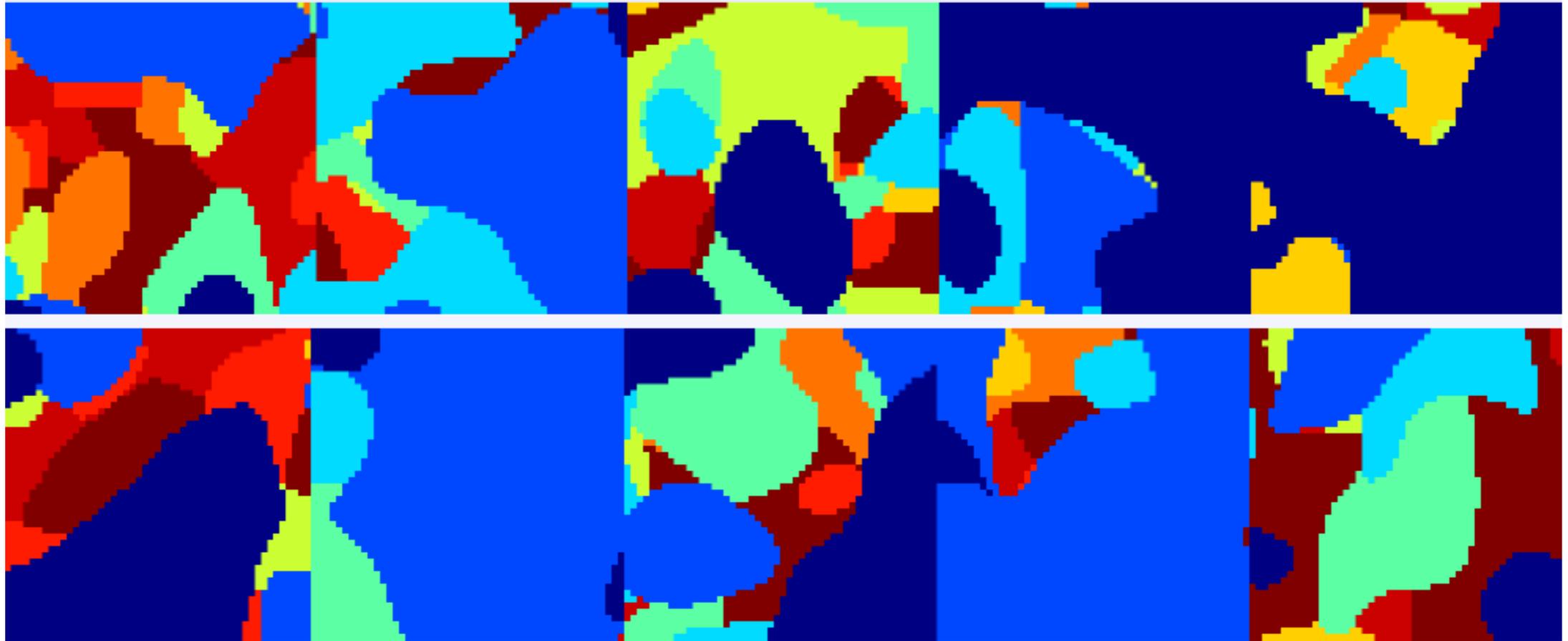


Expected Layer Appearance

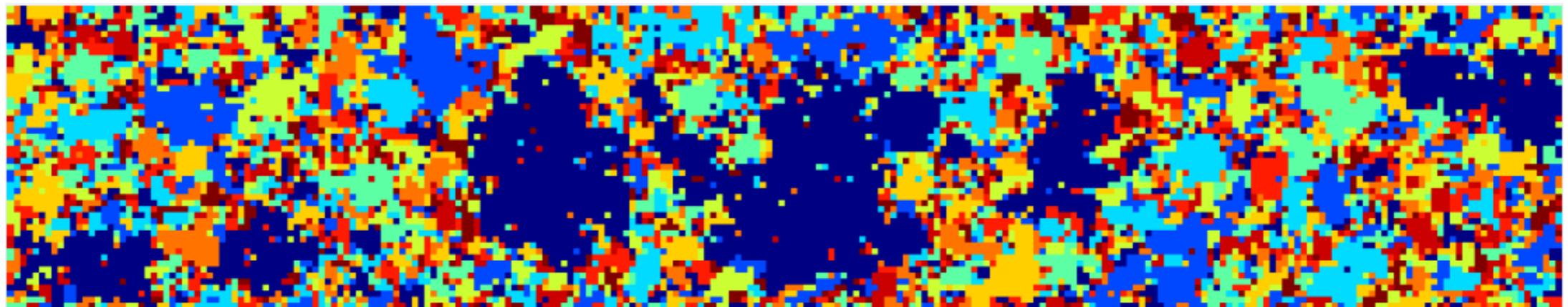


Ghosh & Sudderth, CVPR 2012
 Sudderth & Jordan, NIPS 2008

Generative Samples



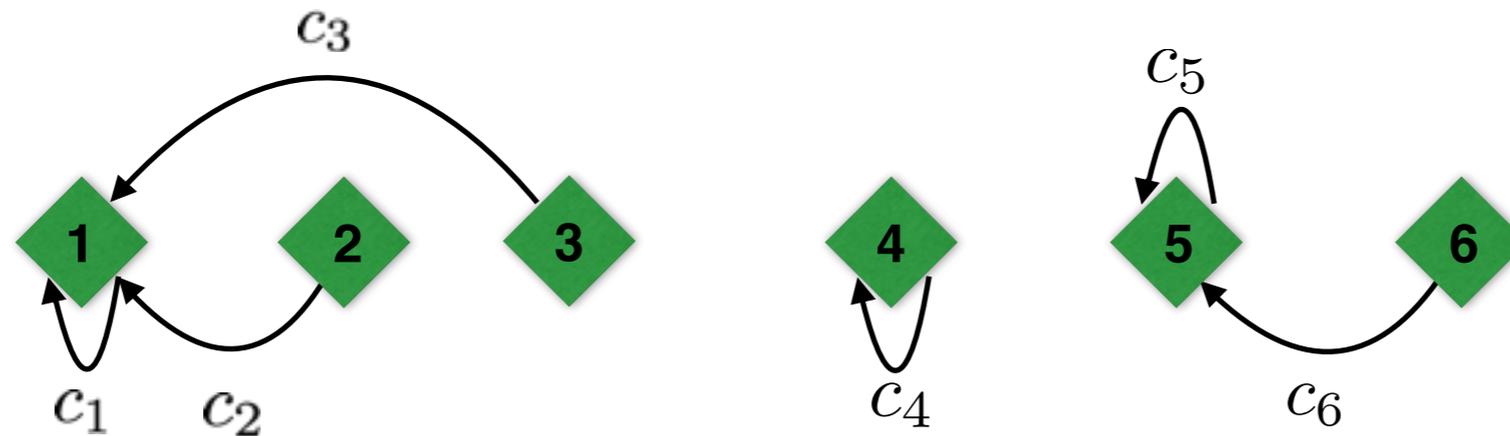
Samples from a Potts Markov Random Field (MRF) model:



Talk Outline

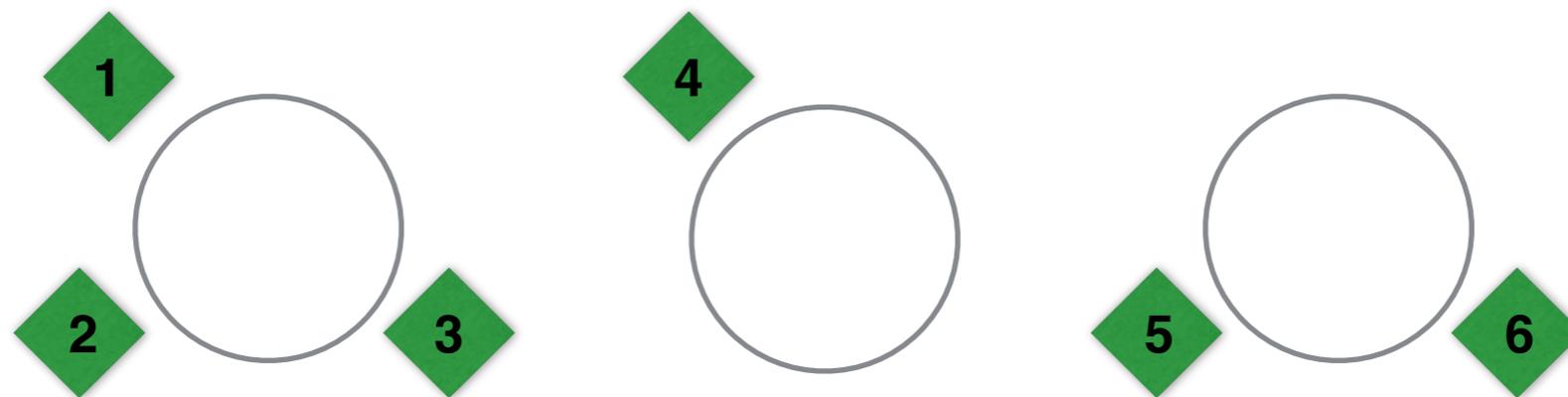
- Distance dependent partitions
 - Parts from articulated 3D objects
- Hierarchical distance dependent partitions
 - Activity discovery from MoCap data
- Learning distance dependent models
 - Image and video segmentation

A distribution over partitions: Chinese Restaurant Process



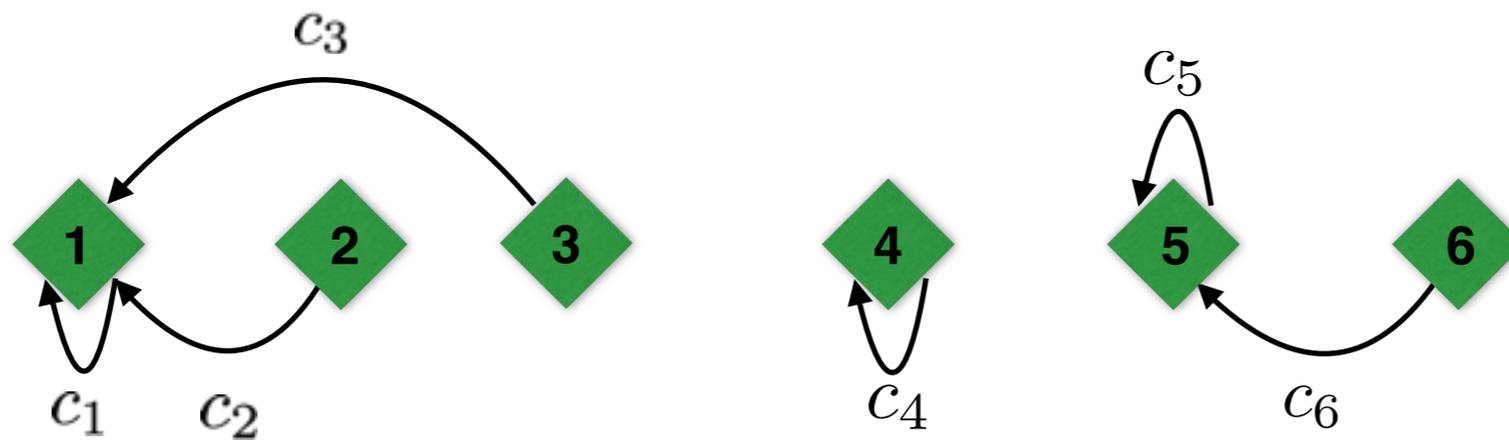
$z(c)$

Customers = Data Instances
Tables = Components



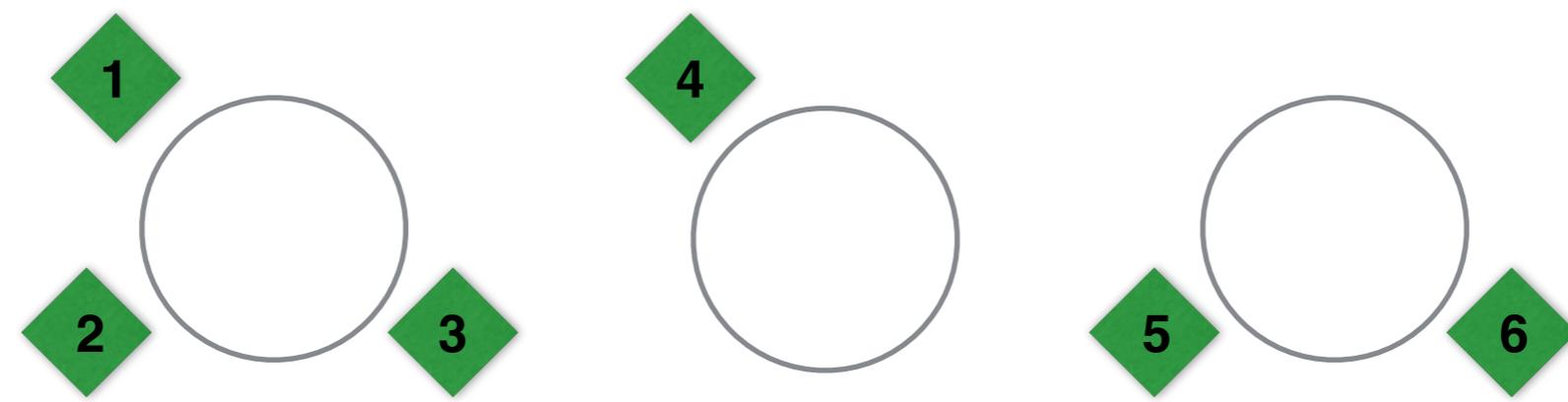
Probability of a customer joining a table $\propto \begin{cases} n_k & \text{if } k \text{ is an existing table} \\ \alpha & \text{if } k \text{ is a new table} \end{cases}$

Distance dependent Chinese Restaurant Process (ddCRP)



$z(c)$

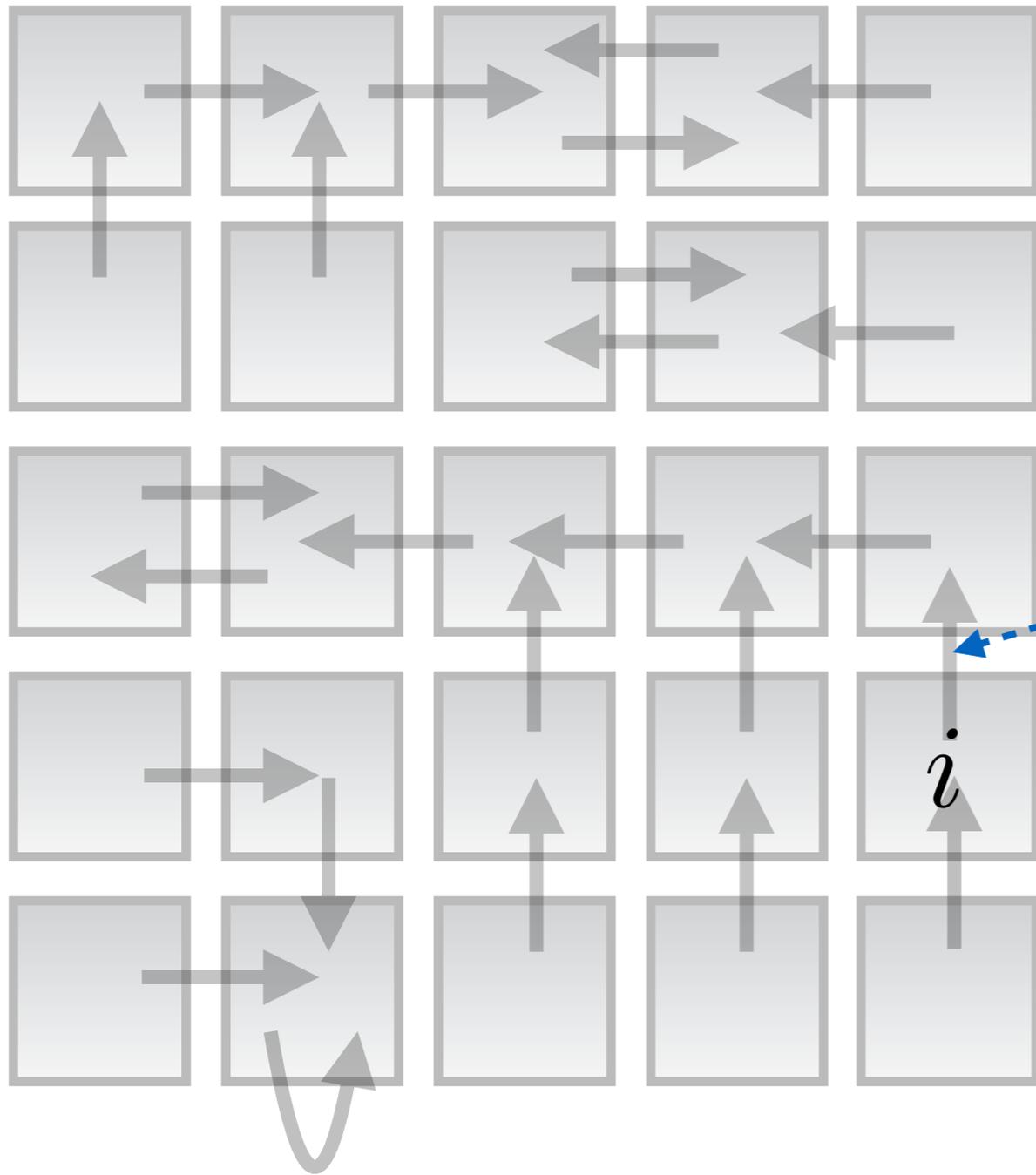
Customers = Data Instances
Tables = Components



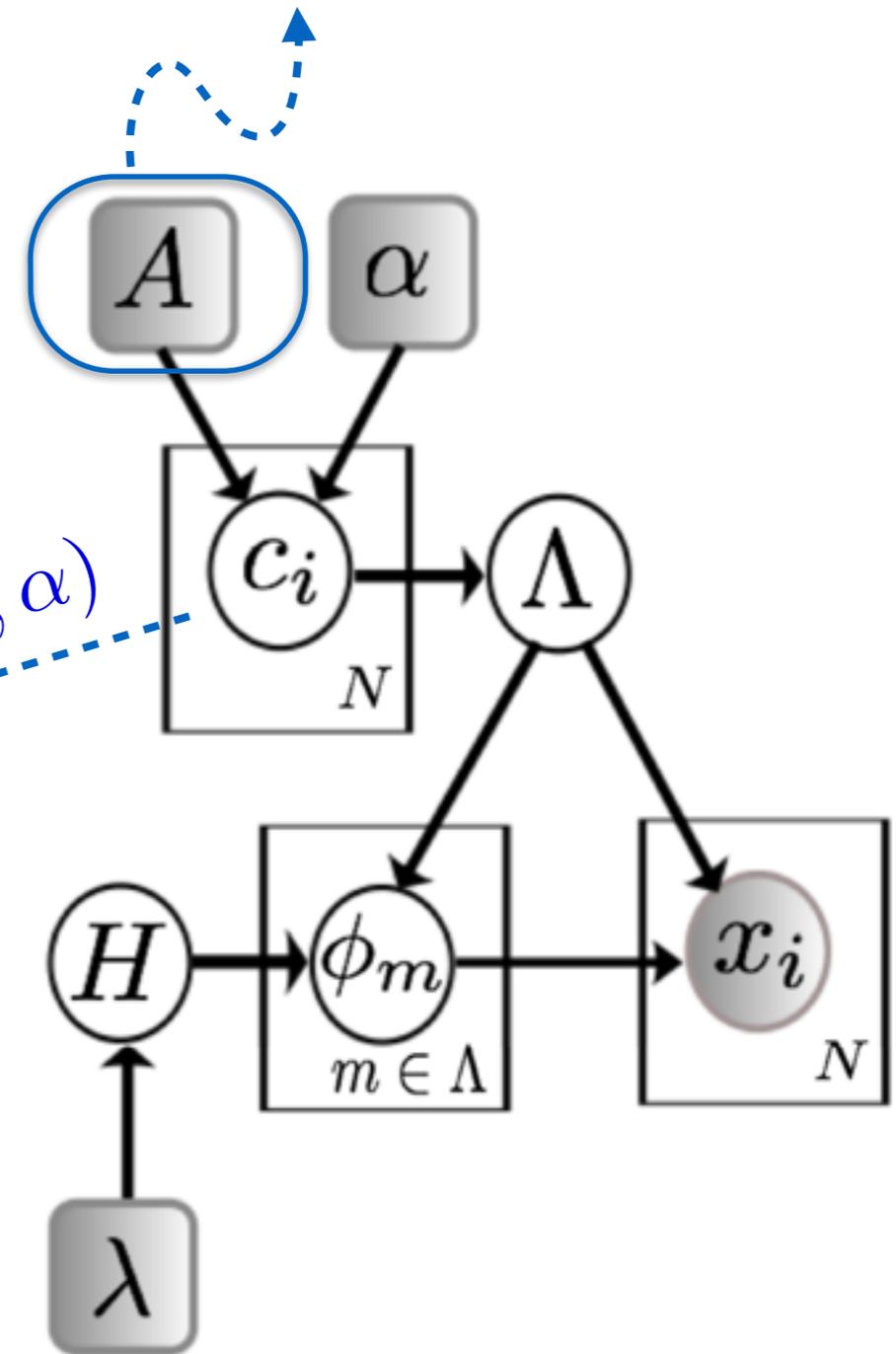
$$p(\mathbf{z}(\mathbf{c}) | \mathbf{A}, \alpha) \propto \prod_{n=1}^N \begin{cases} A_{mn} p(c_n | A, \alpha) & \text{if } m \neq n, \\ \alpha & \text{if } m = n. \end{cases}$$

Models for heterogeneous data

Captures dependencies

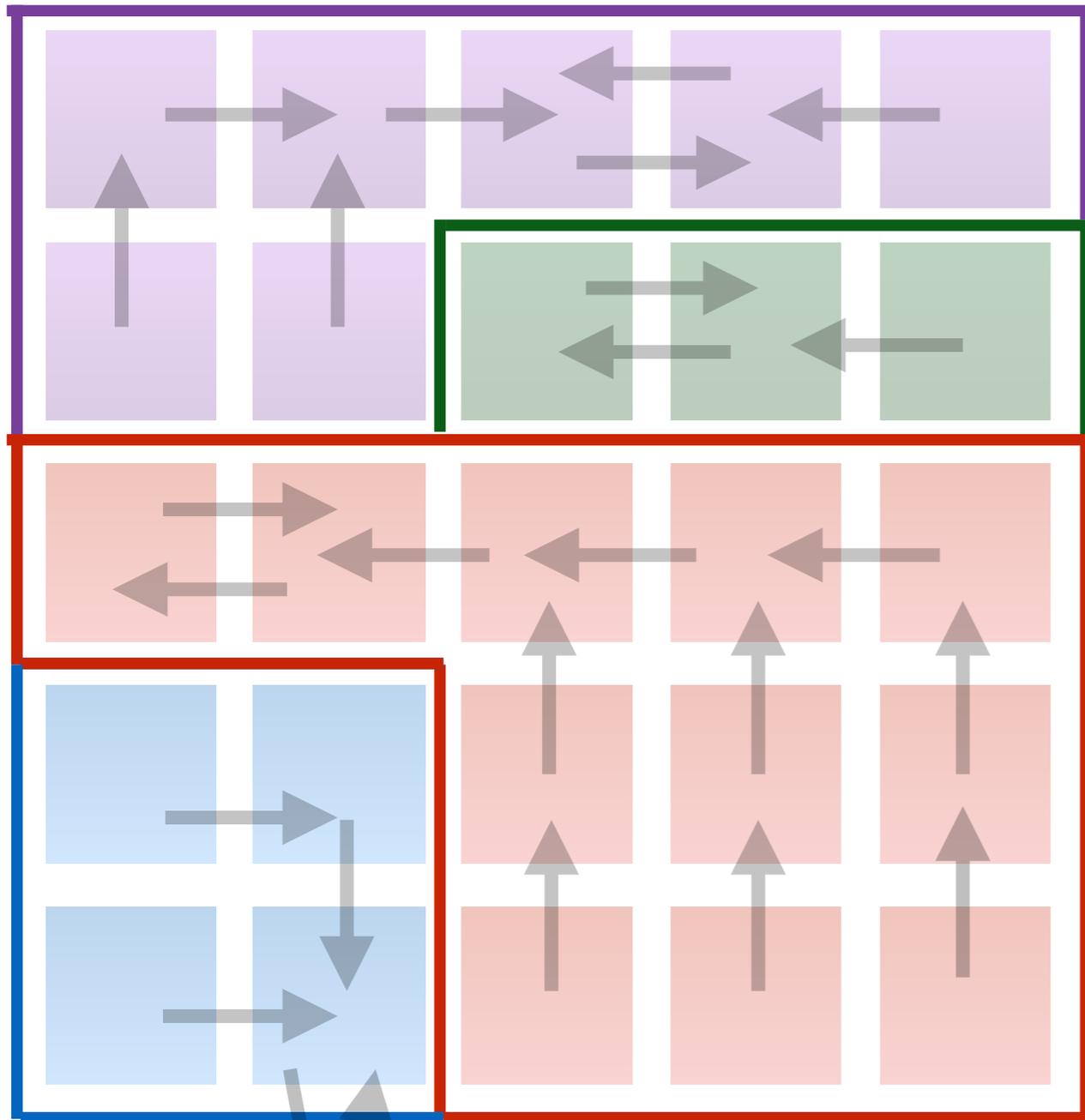


$$c_i \sim p(c_i | A, \alpha)$$

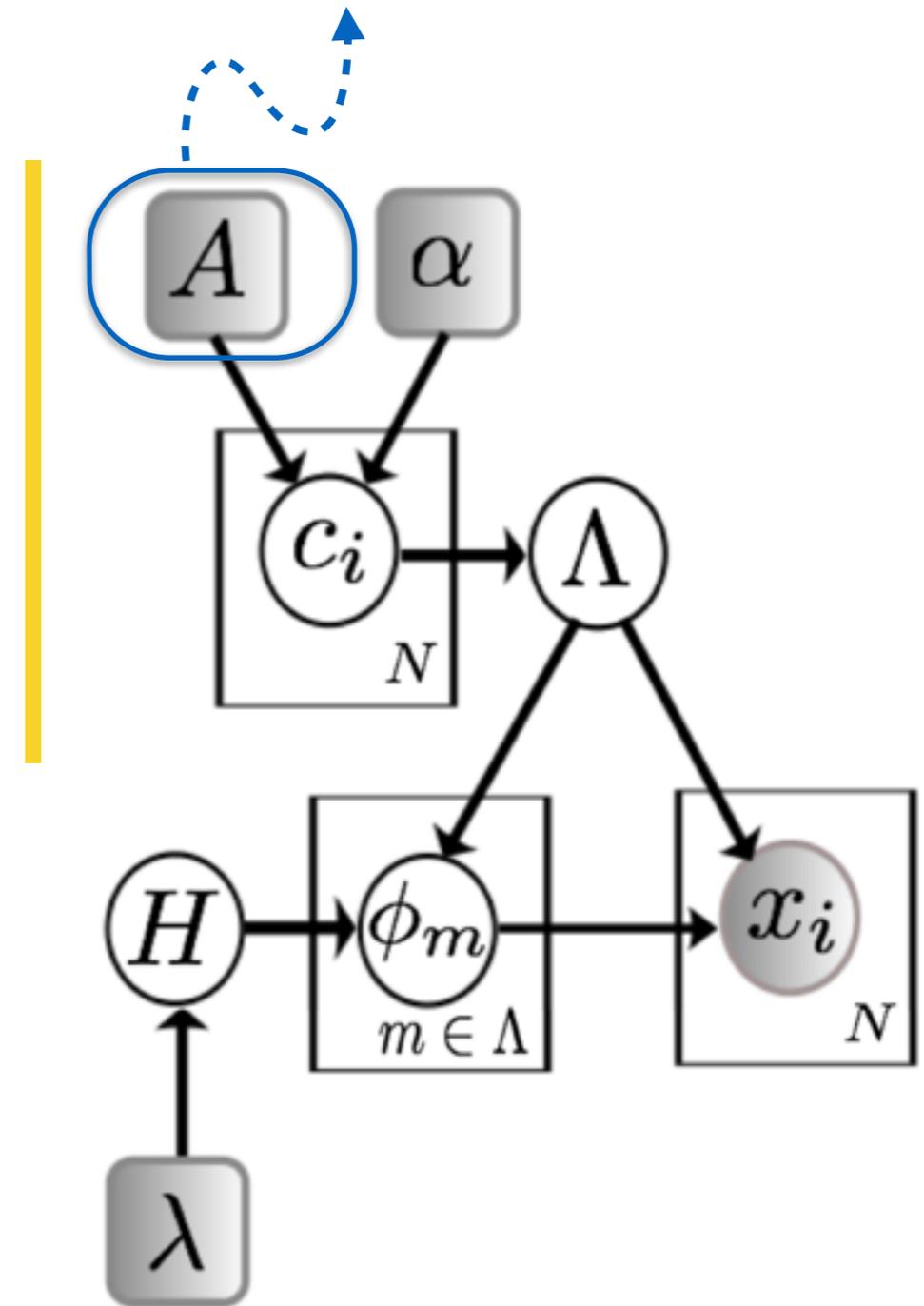


Models for heterogeneous data

Captures dependencies

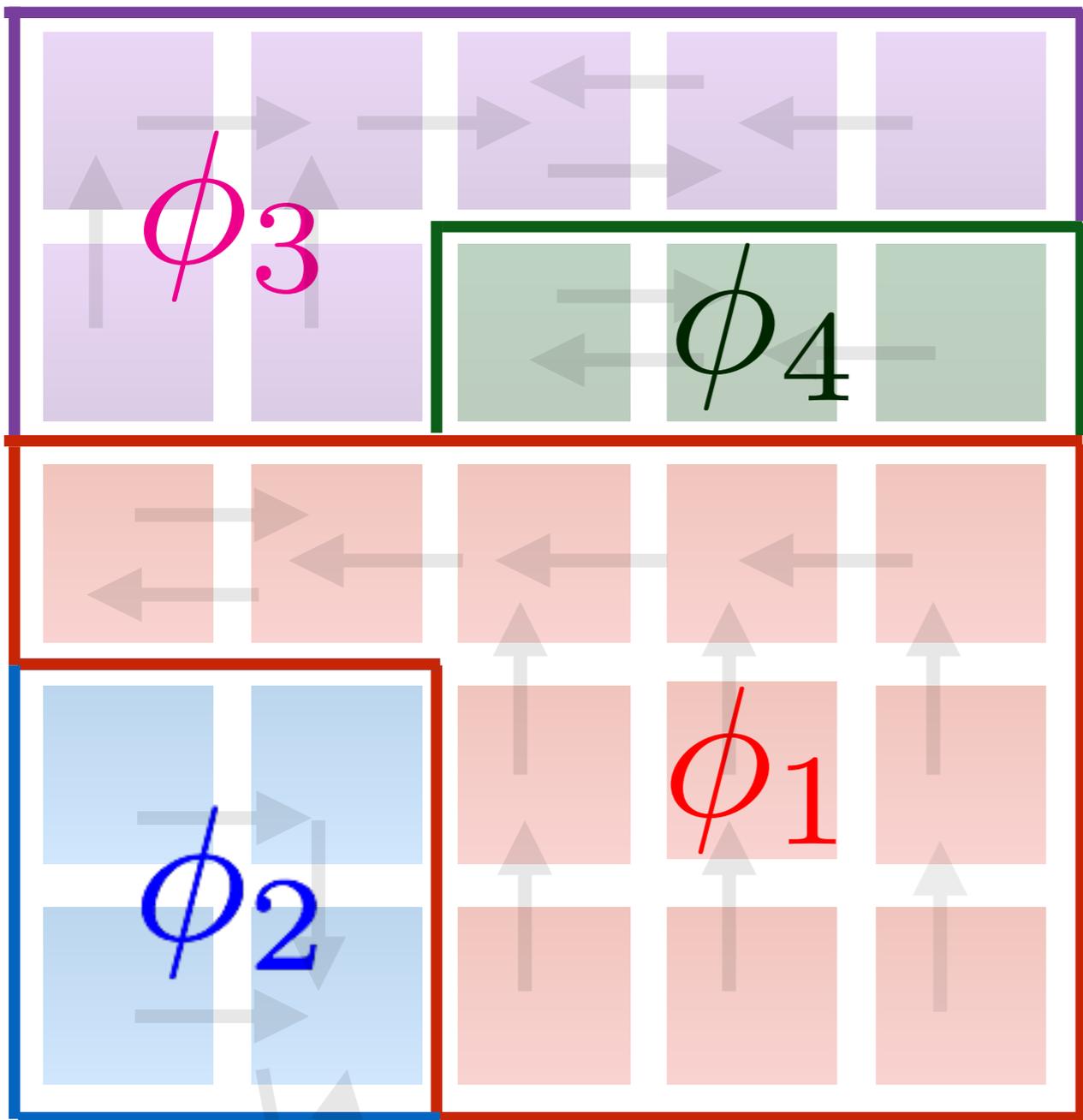


Λ

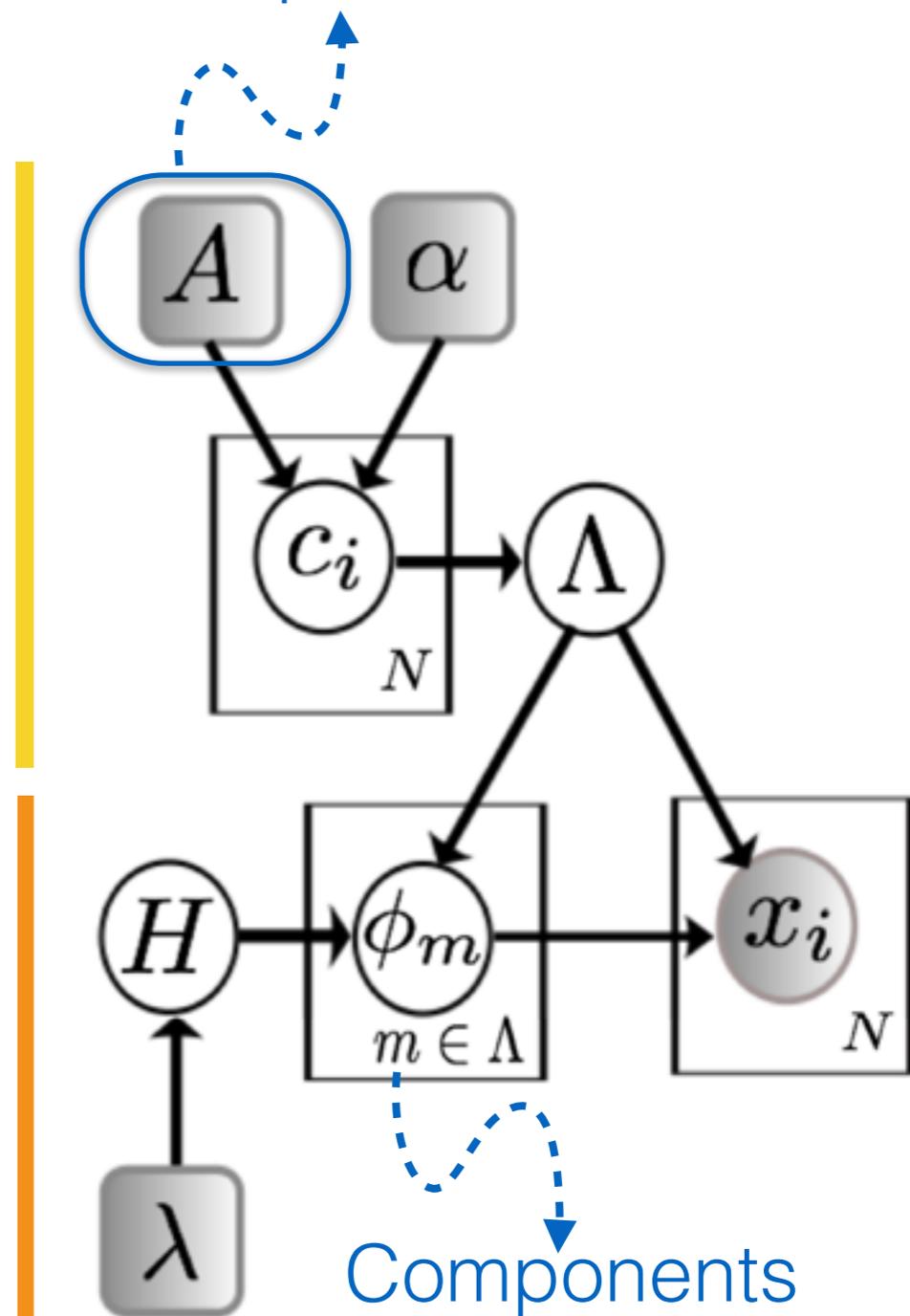


Models for heterogeneous data

Captures dependencies



Λ



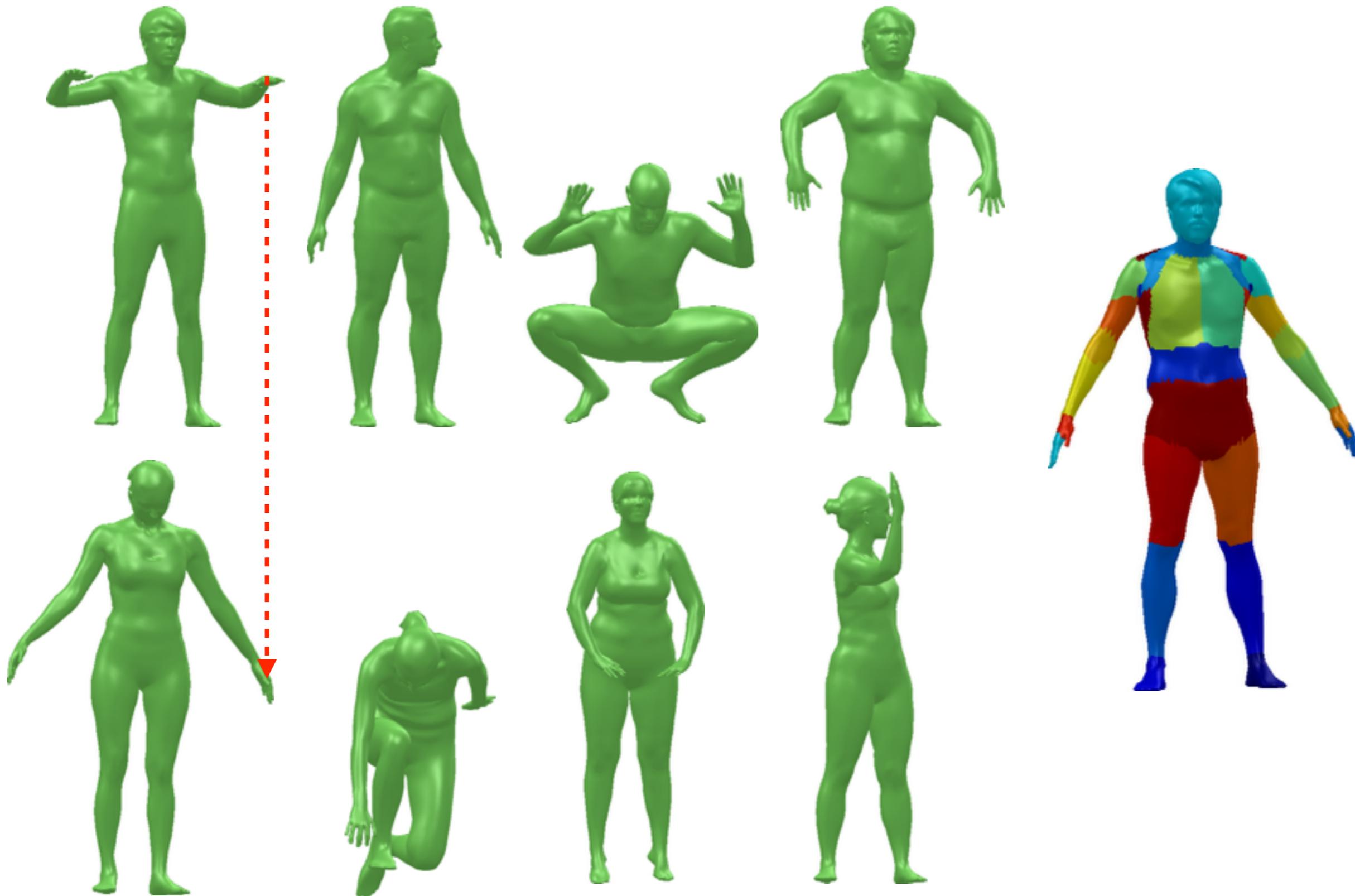
$$\phi_m \sim H(\lambda), \forall m \in \Lambda$$

$$x_i \sim \phi_m, \forall i \in m$$

Talk Outline

- Distance dependent partitions
 - Parts from articulated 3D objects
- Hierarchical distance dependent partitions
 - Activity discovery from MoCap data
- Learning distance dependent models
 - Image and video segmentation

Parts from Deformations

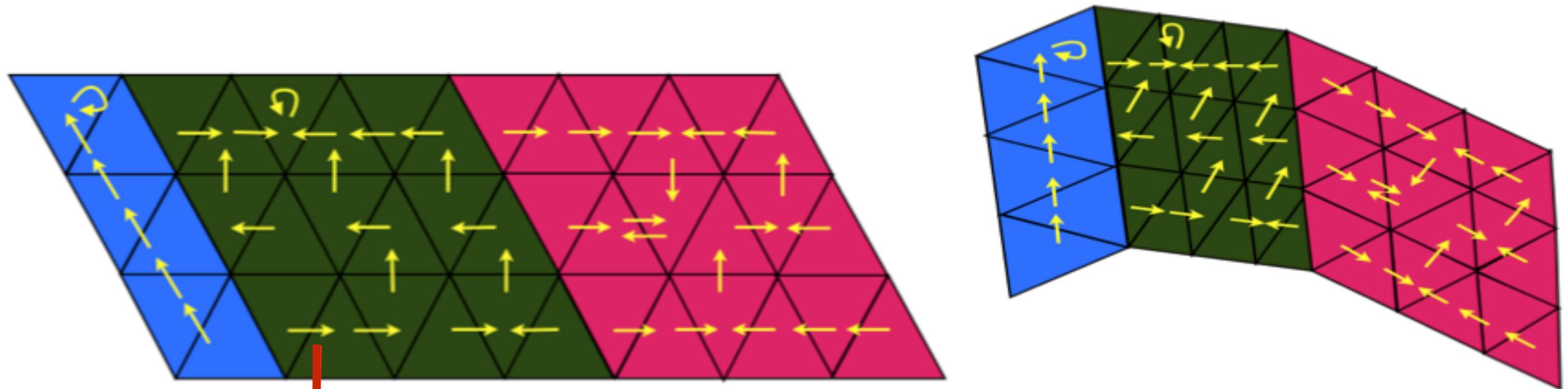


Discovering Parts from Deformations: Big Picture



- **Cluster:** Mesh faces.
- **Prior:** over the space of plausible mesh partitions.
- **Likelihood:** Given segmentation into parts, model how multiple bodies deform across many poses.
- **Posterior:** Explored through MCMC.

ddCRP Prior over Mesh Partitions



$$p(c_m = n \mid A, \alpha) \propto \begin{cases} A_{mn} & m \neq n, \\ \alpha & m = n. \end{cases}$$

- Mesh faces are only allowed to link to neighboring faces

$$A_{mn} = \mathbf{1}[d_{mn} \leq 1]$$

Prior over plausible partitions



$$p(Z_1) > 0$$

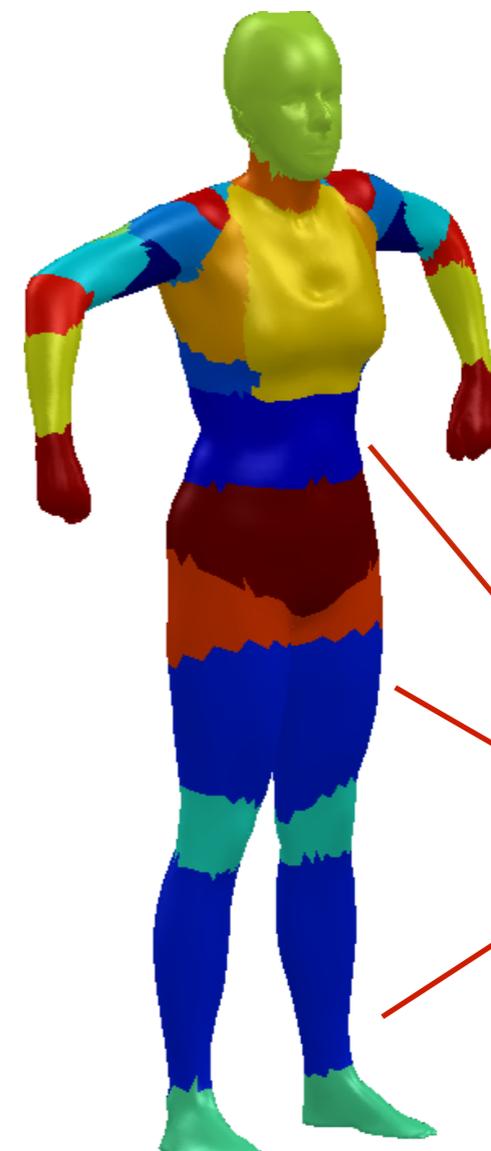
Desirable

Noise



$$p(Z_2) = 0$$

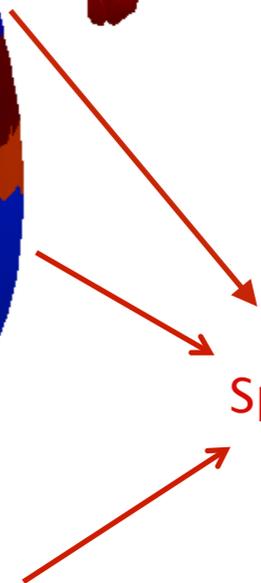
Avoid: Noisy Parts



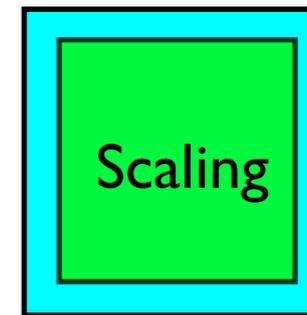
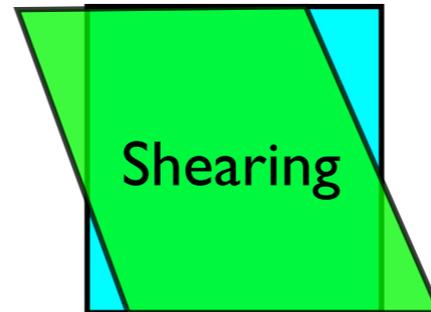
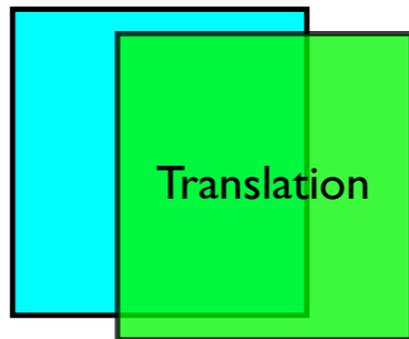
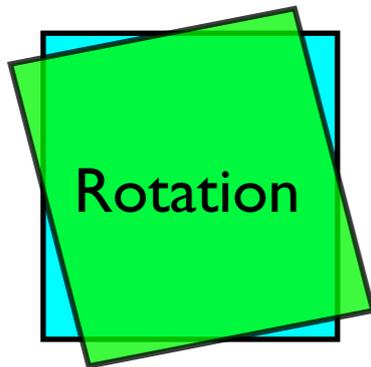
$$p(Z_3) = 0$$

Avoid: Disconnected Parts

Split limbs



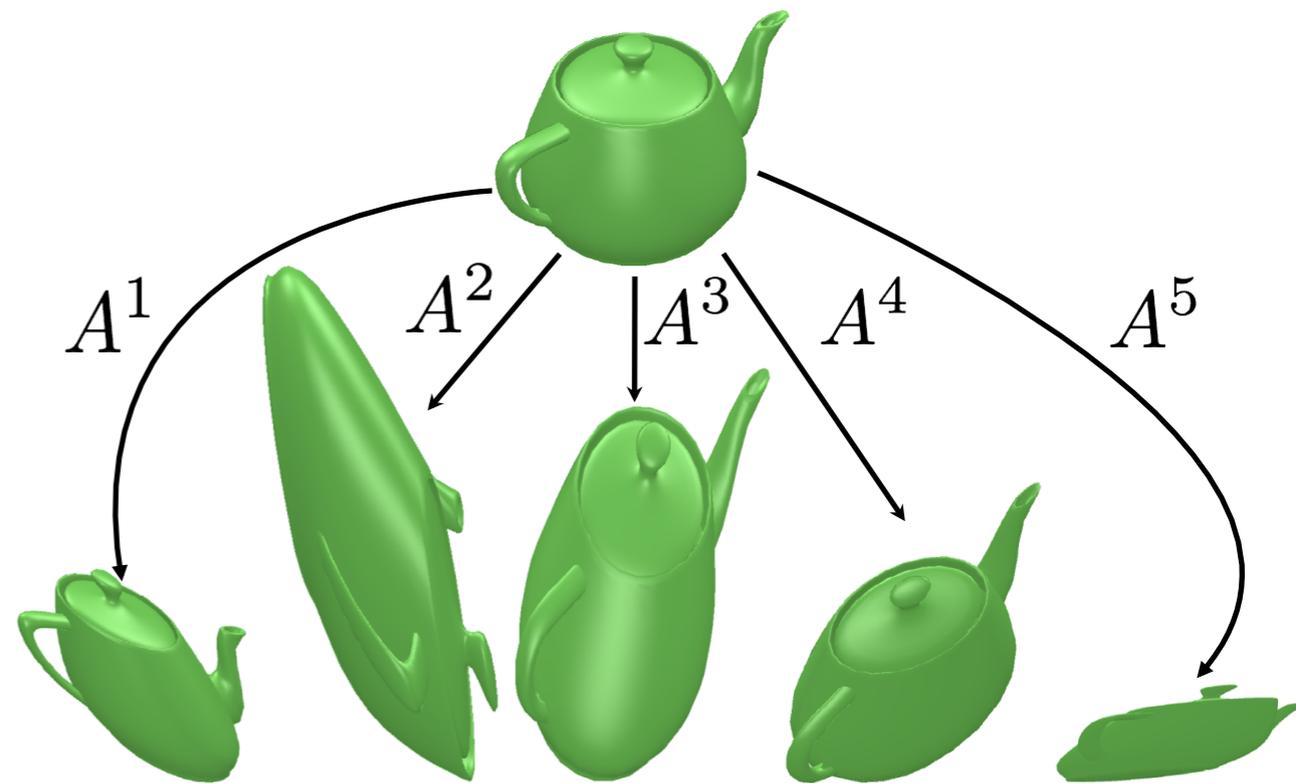
Modeling Part Deformations



Matrix Normal Inverse Wishart:

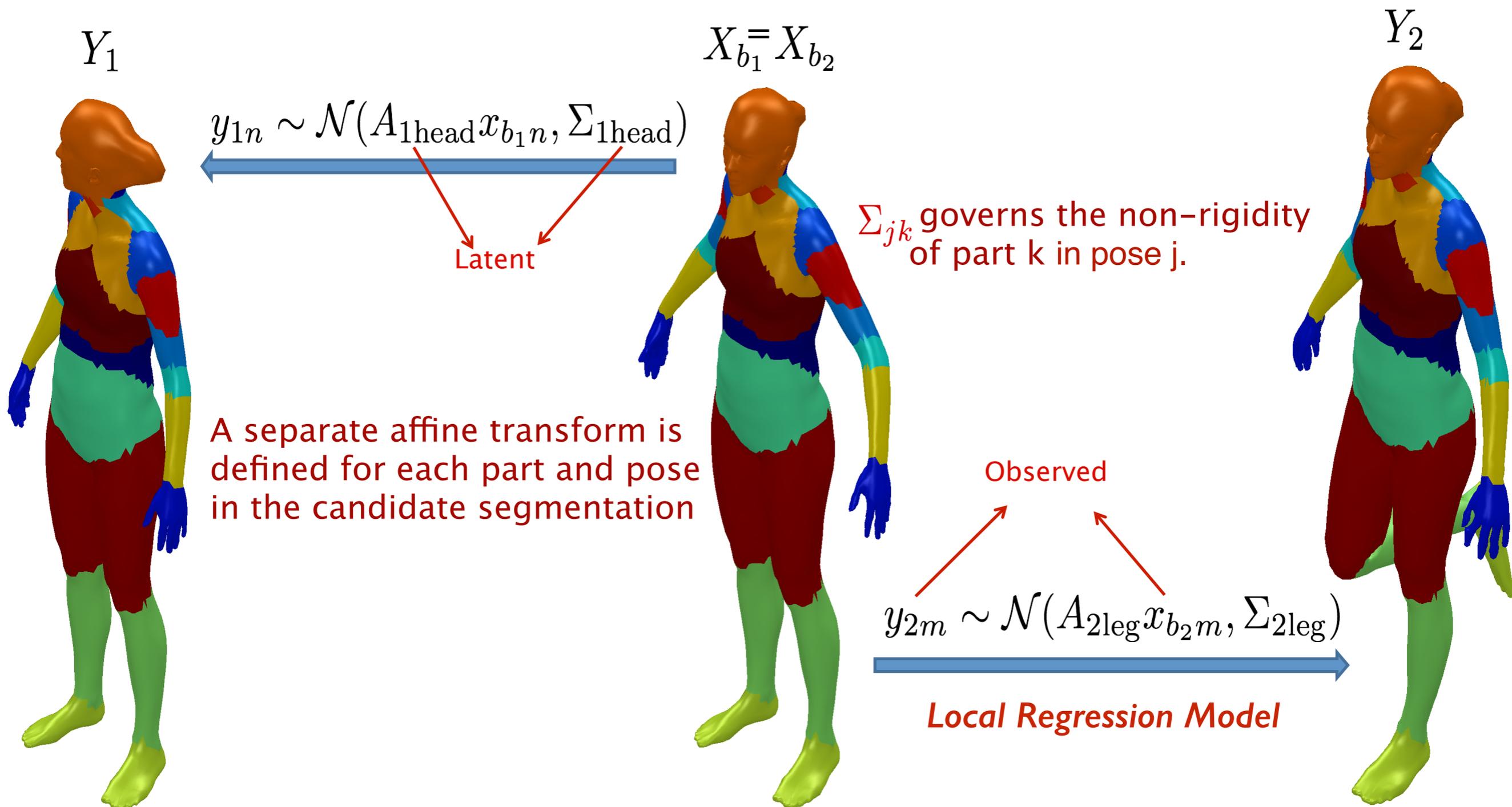
$$\Sigma \sim \mathcal{IW}(n_0, S_0)$$
$$A \mid \Sigma \sim \mathcal{MN}(M, \Sigma, K)$$

where $A \in \mathbb{R}^{3 \times 4}$ is an **affine** transformation.



$$A^1 \dots A^5 \sim \mathcal{MNIW}(\cdot)$$

Generative Affine Likelihoods



Marginal Affine Likelihoods

For each part and pose combination analytically marginalize over *all possible* affine transformations

$$p(Y_{jk} | X_{jk}) = \int p(Y_{jk}, A_{jk}, \Sigma_{jk} | X_{jk}) dA_{jk}, d\Sigma_{jk}$$

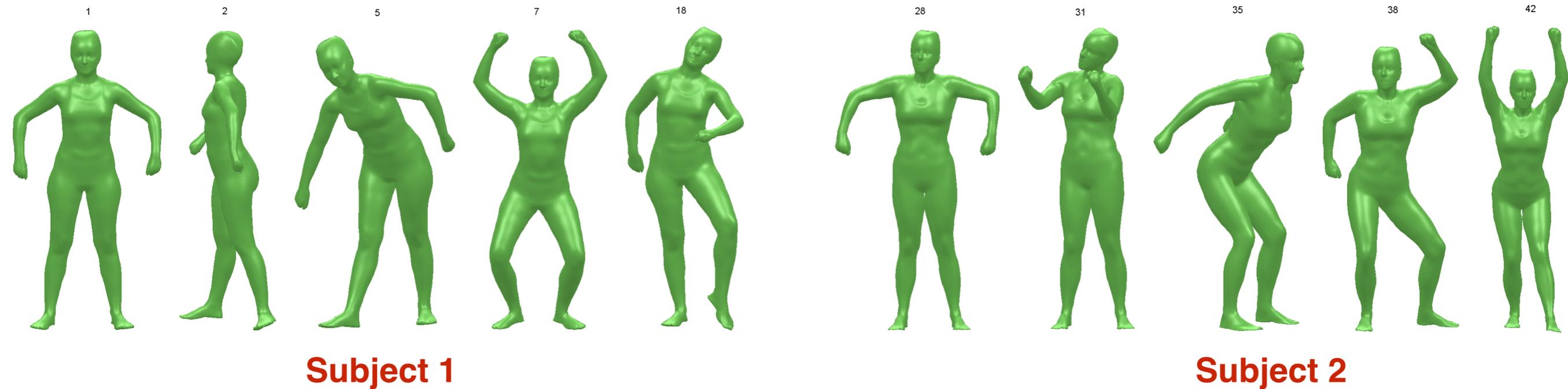
Marginal Likelihood

Bayesian Model Selection:

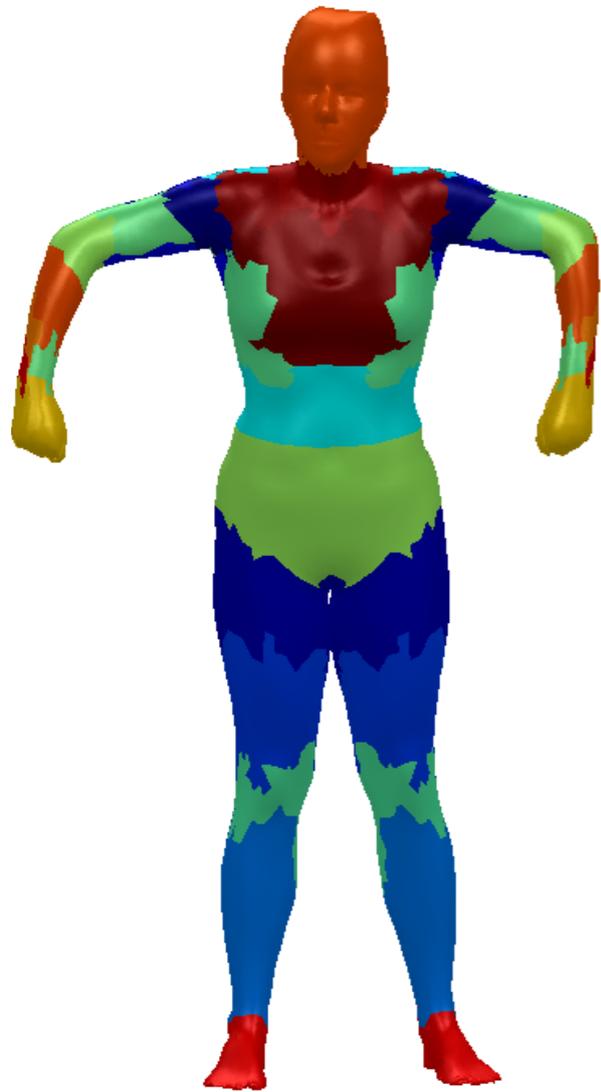
- Improper merges have low marginal likelihoods
- Improper splits are “suspicious coincidences” and end up with lower marginal likelihoods

Human Bodies in Motion

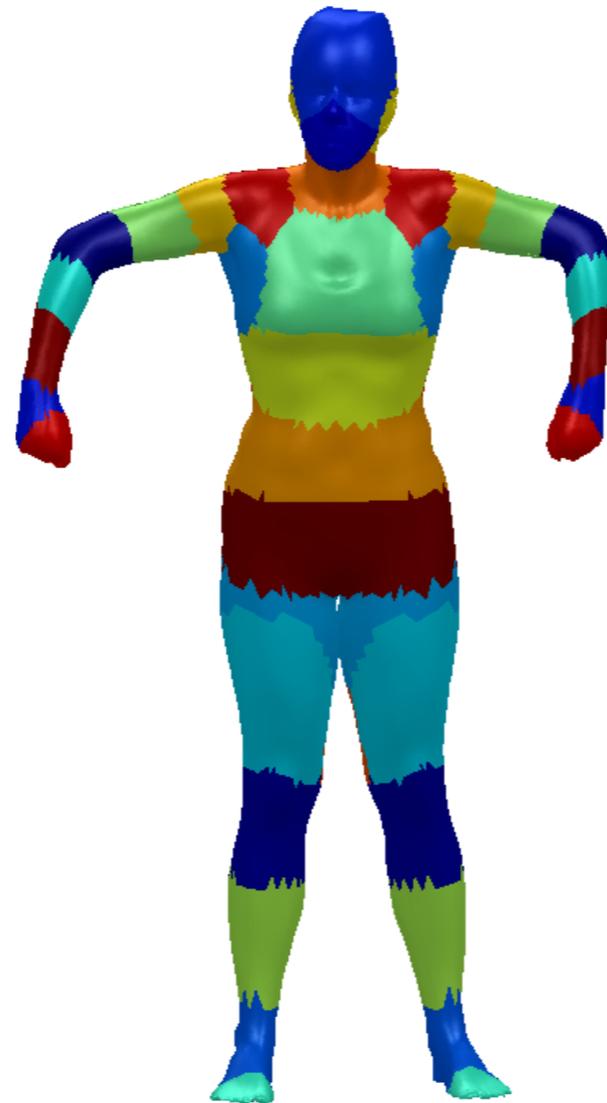
- **56** Aligned scans from **two** human subjects
- Wide variability in poses, limited variability in body shapes



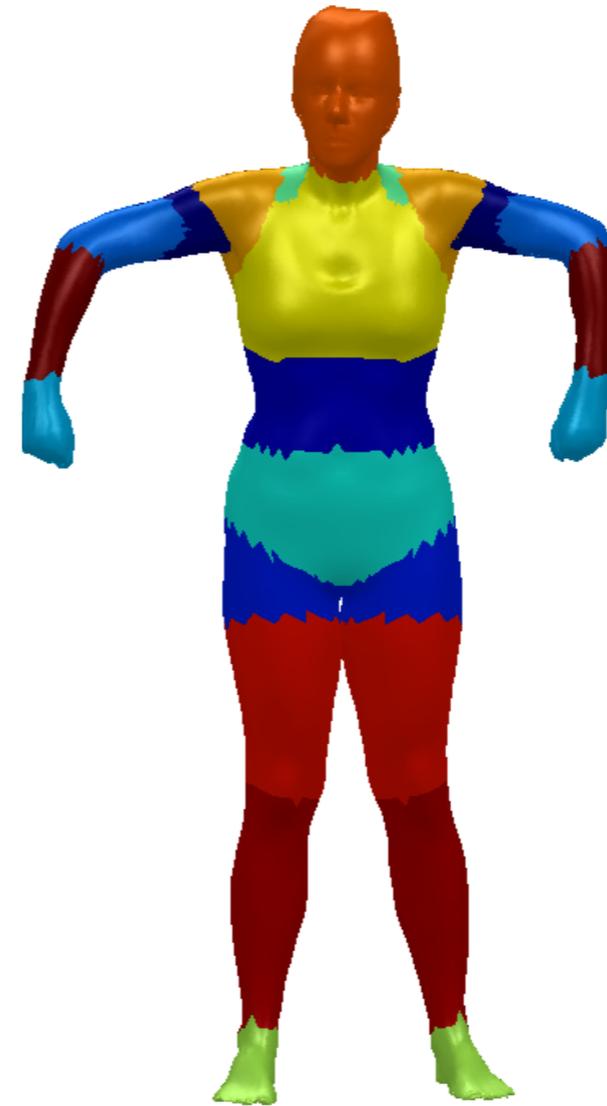
Visual Comparisons



Agglomerative



Spectral Clustering
Liu & Zhang, 2004

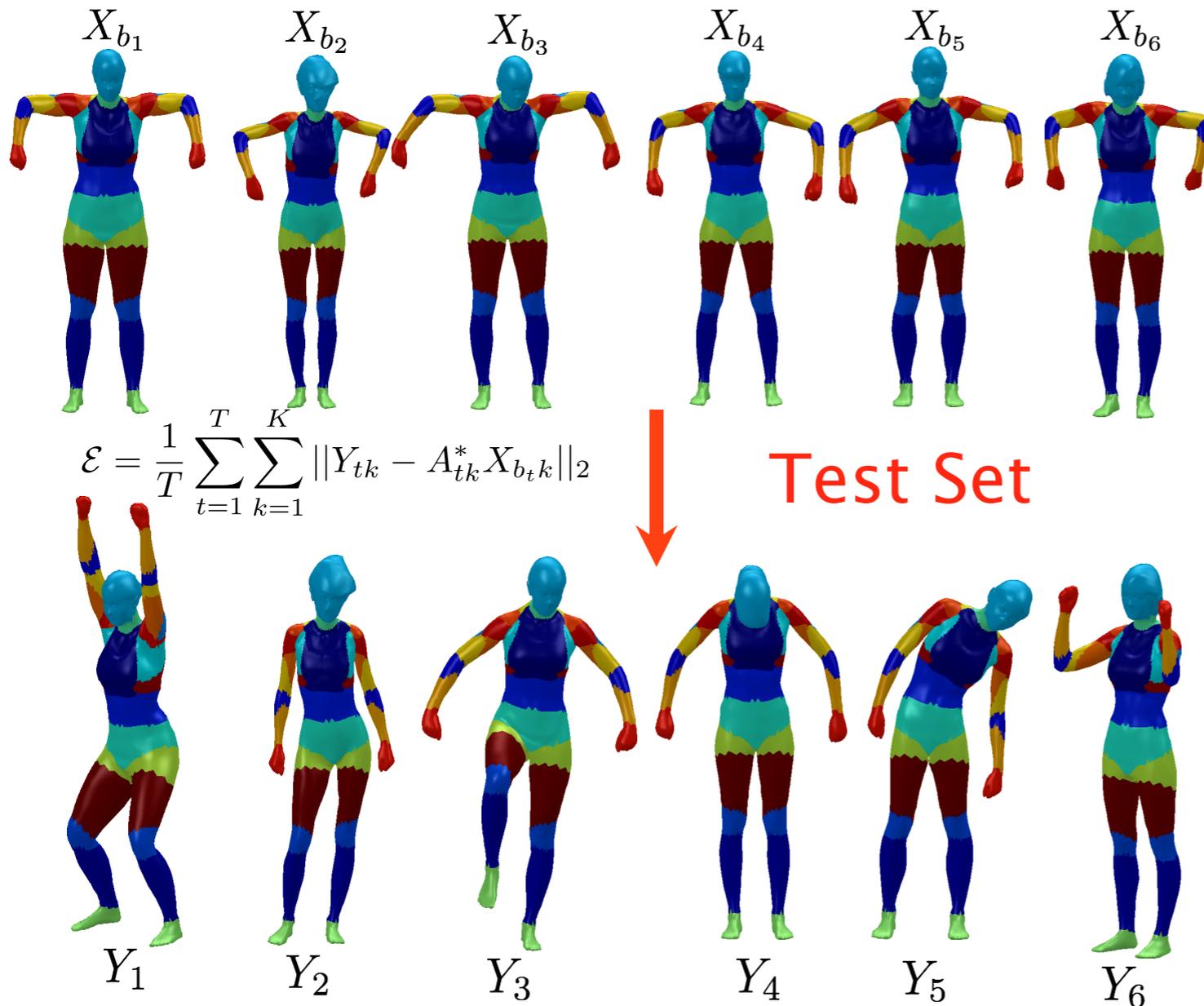


CRP



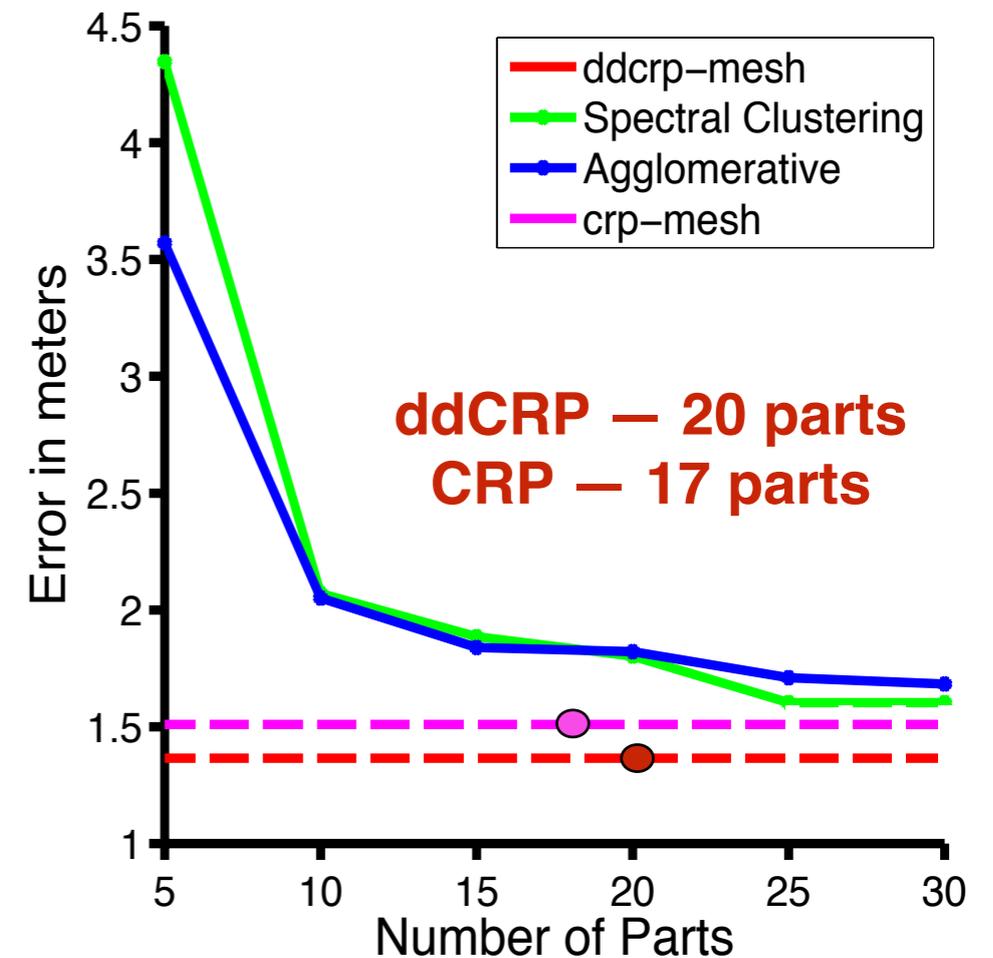
ddCRP

Quantitative Evaluation

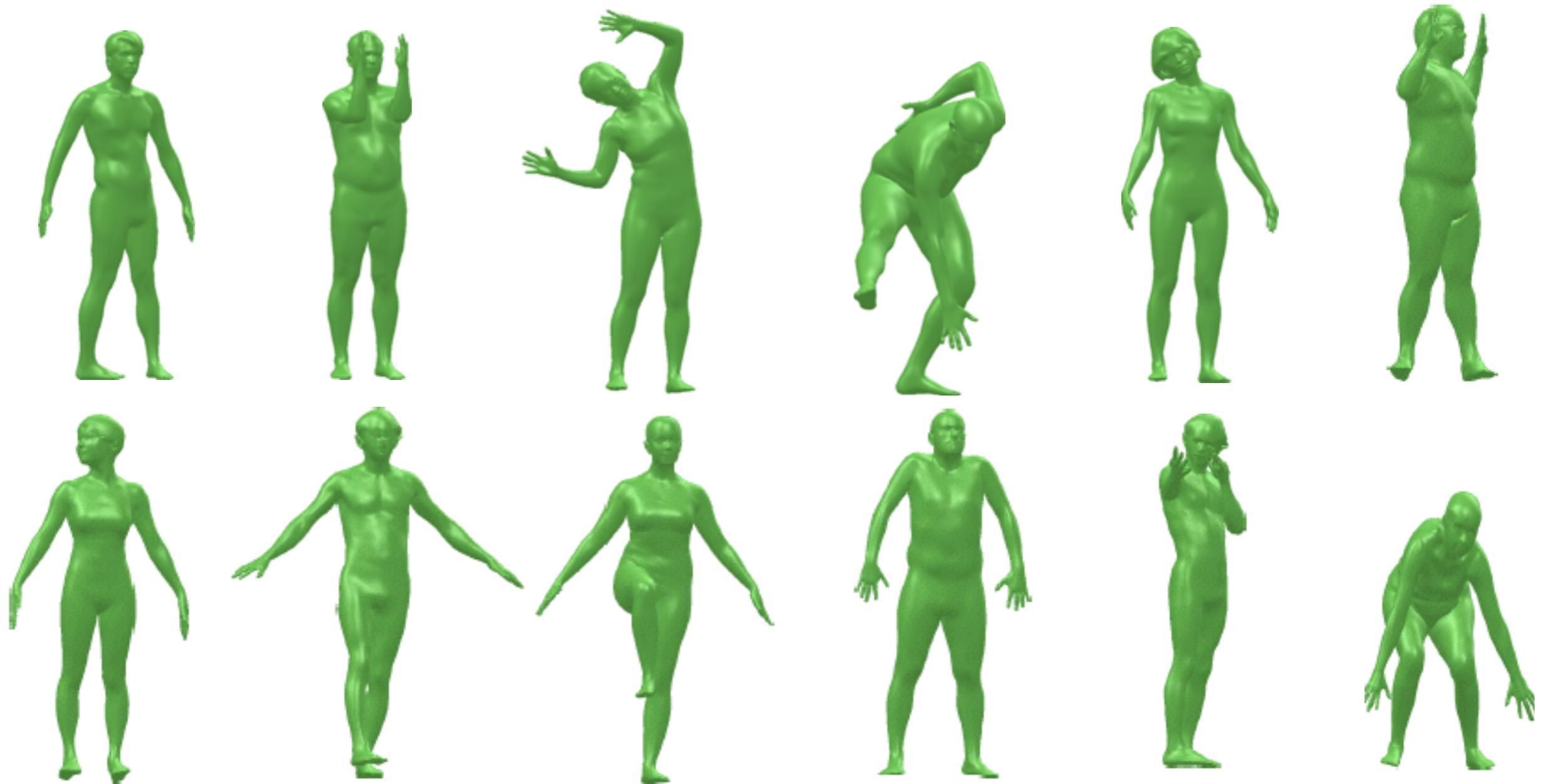


12 Test meshes from six subjects.

Measure error in predicted motion for the candidate segmentations



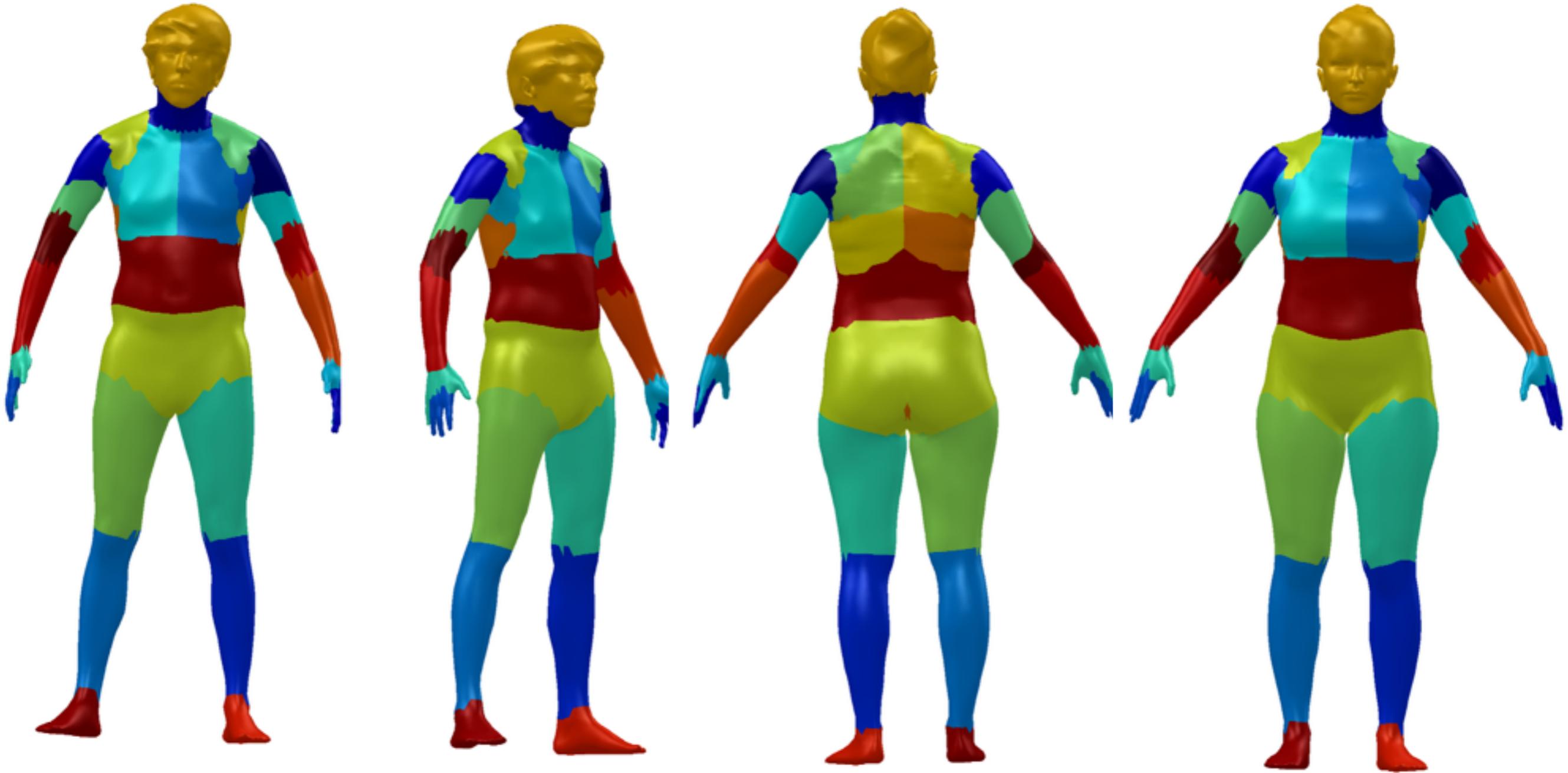
Large Scale Studies



1732 meshes, 78 subjects, ~22,000 mesh faces

Wide variability in both body shapes and poses.

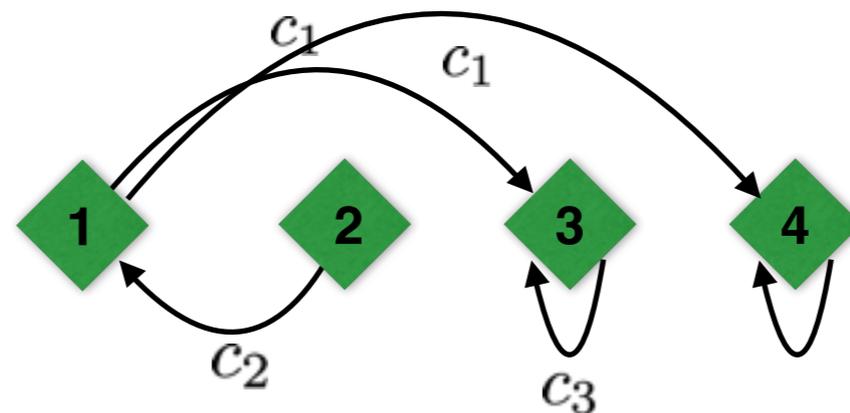
Segmented Bodies



Computer generated meshes



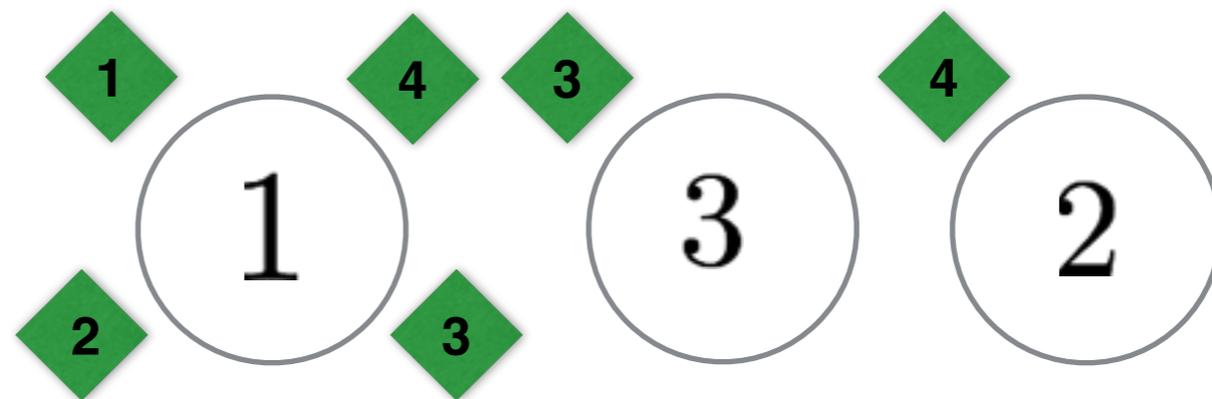
Inference through Gibbs Sampling



*Customers = Mesh Faces
Tables = Object Parts*

Collapsed Sampler:
Only need to sample links,
other random variables are
analytically marginalized out.

$z(c)$



*Table structure
↕
Segmentation*

Local changes in the link structure lead to large changes in the partition structure

Talk Outline

- Distance dependent partitions
 - Parts from articulated 3D objects
- Hierarchical distance dependent partitions
 - Activity discovery from MoCap data
- Learning distance dependent models
 - Image and video segmentation

Hierarchical Distance Dependent Partitions

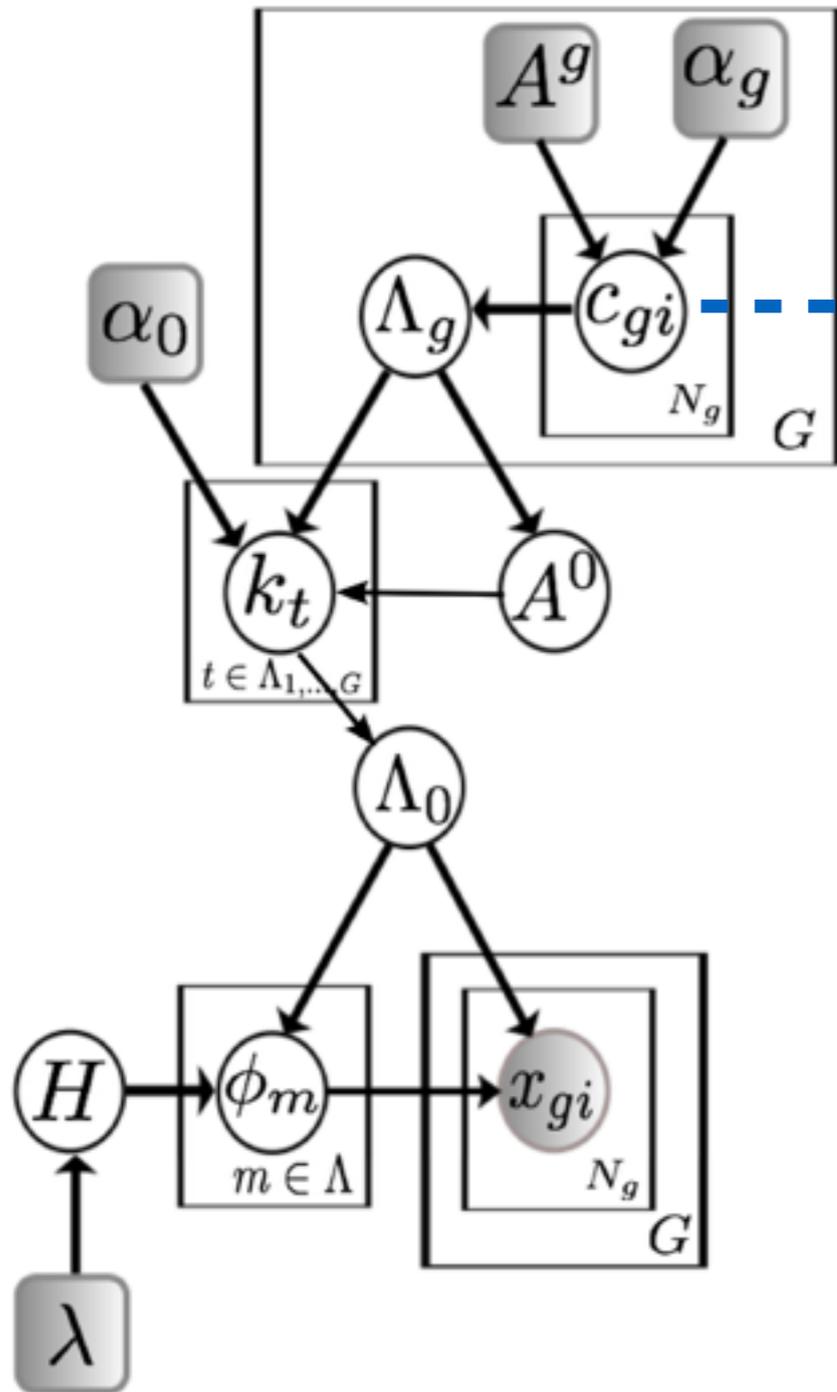
Doctors have been confounded by the divergent paths of Ebola patients whose cases appeared similar at first. They have been especially baffled by the “light bulb” phenomenon ... genes

Doctors have been confounded by the divergent paths of Ebola patients whose cases appeared similar at first. They have been especially baffled by the “light bulb” phenomenon ... genes

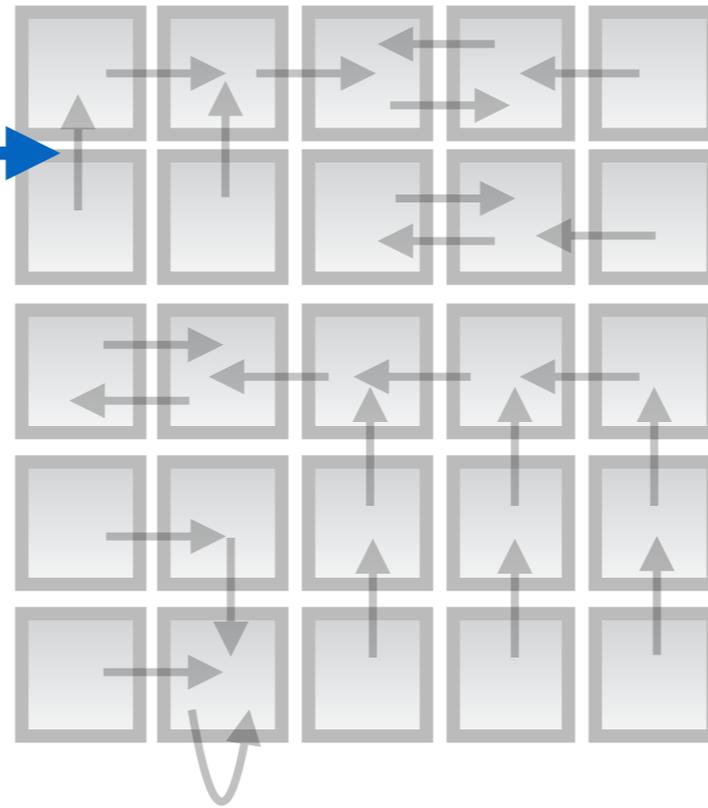


Model affinities between both data points and *latent* clusters.

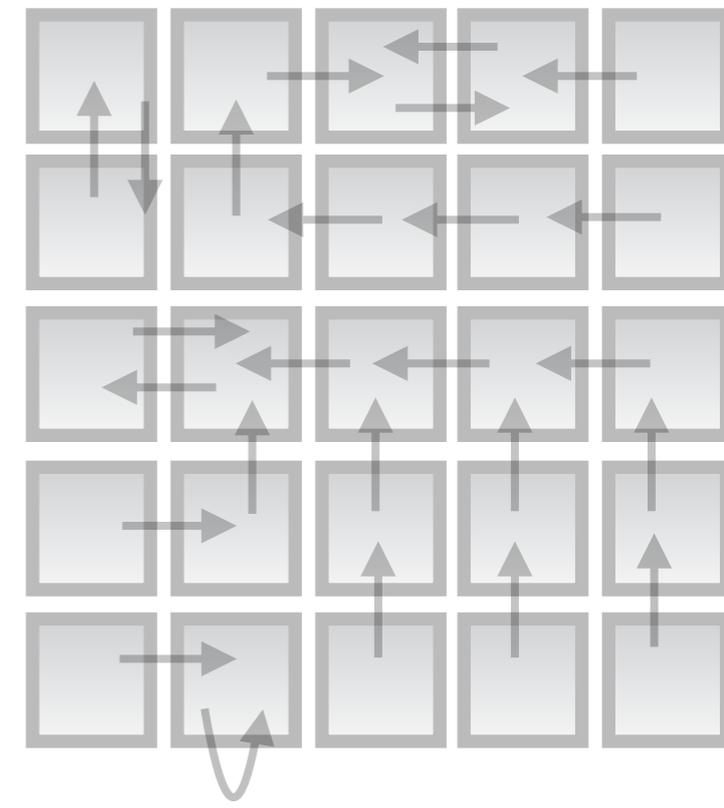
Hierarchical ddCRP



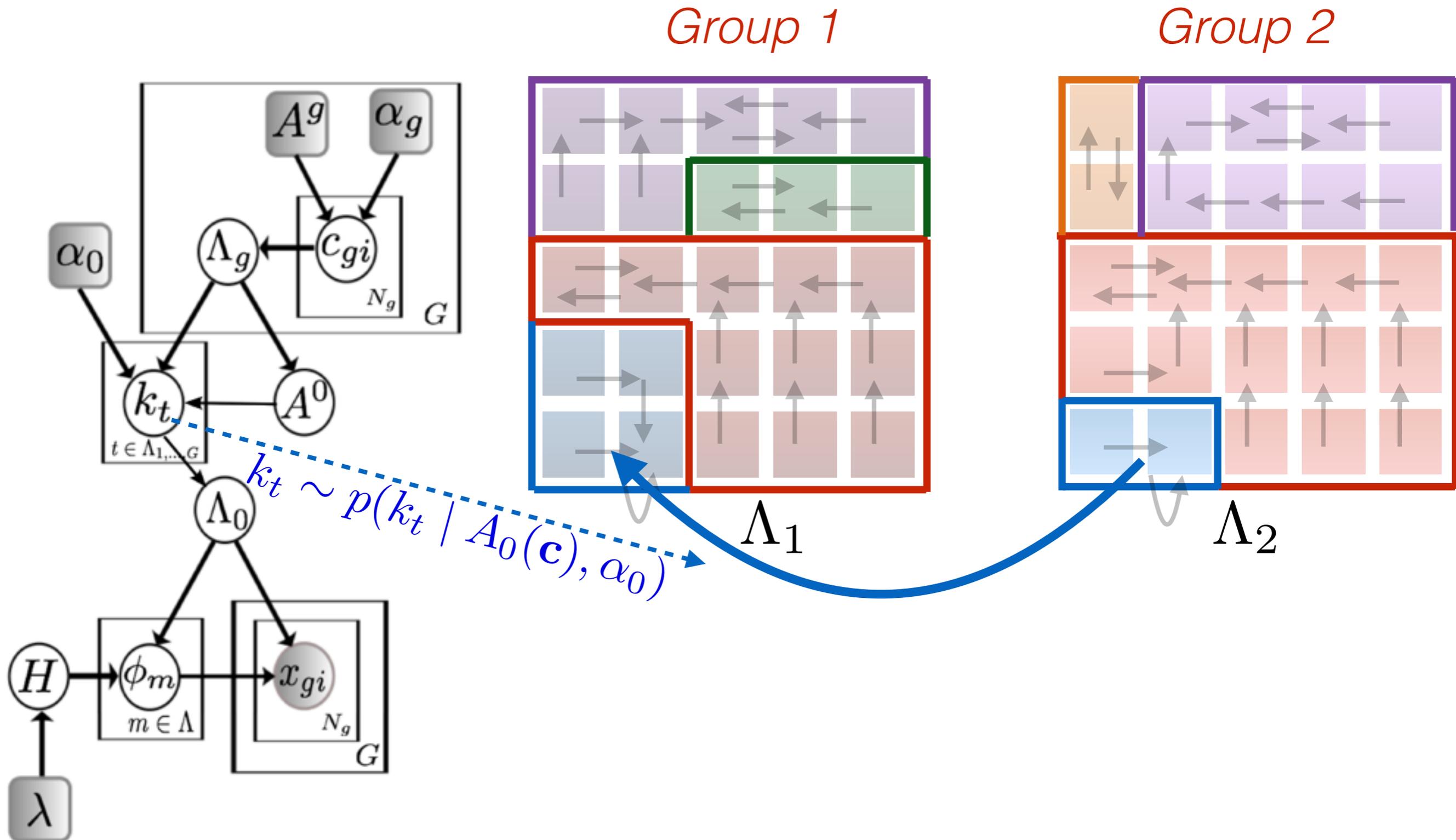
Group 1



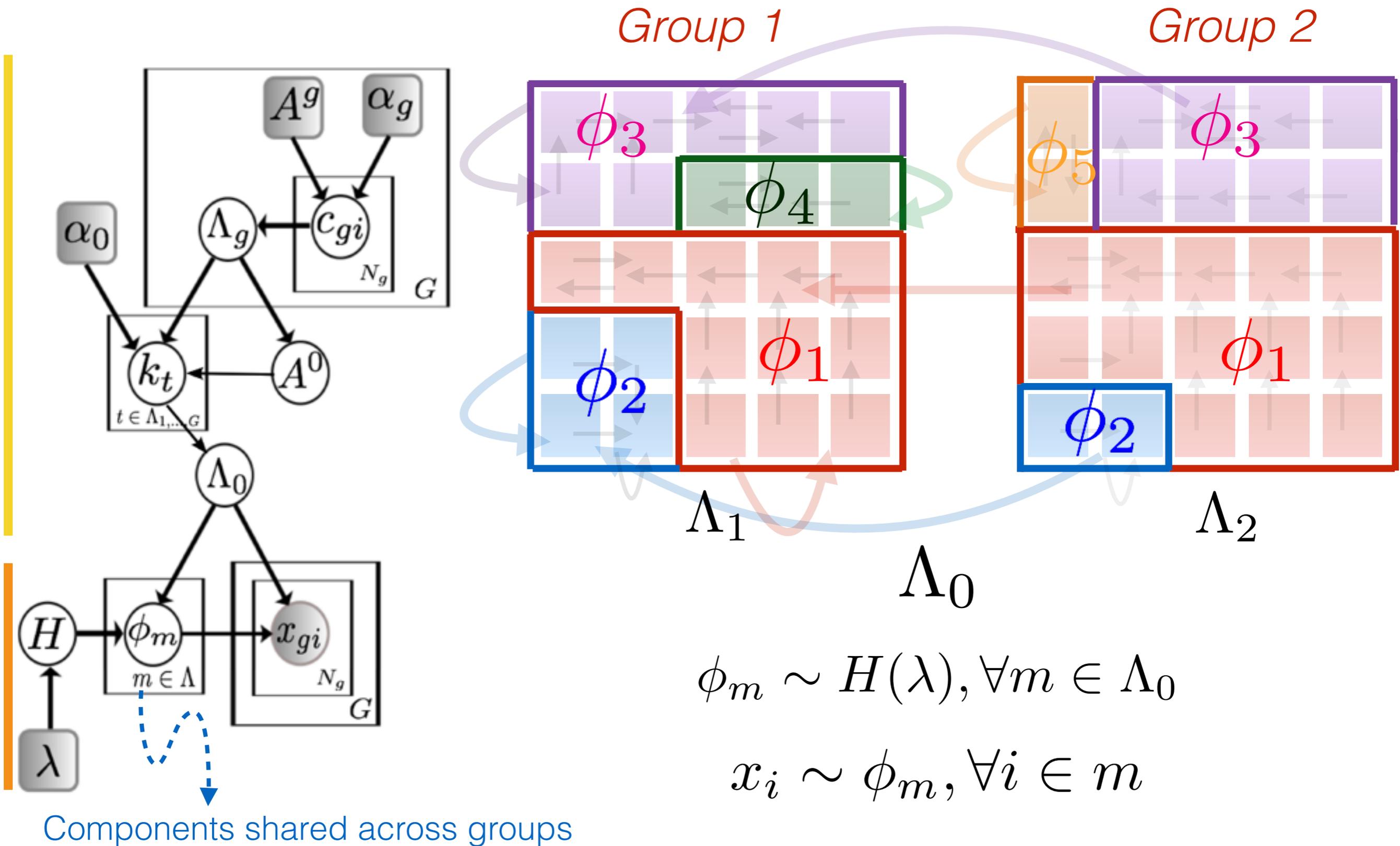
Group 2



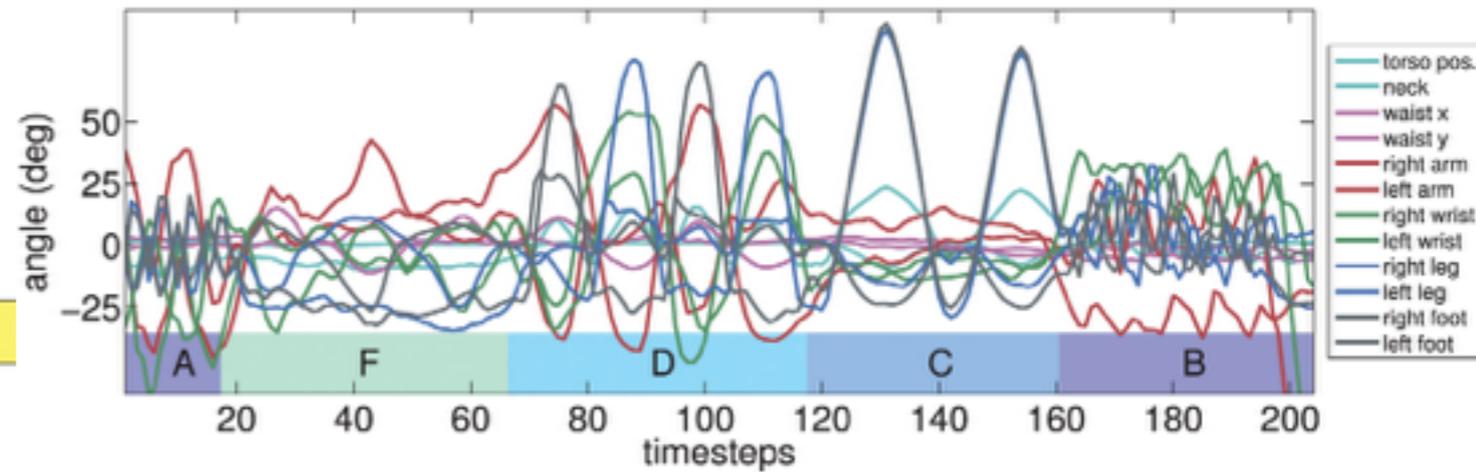
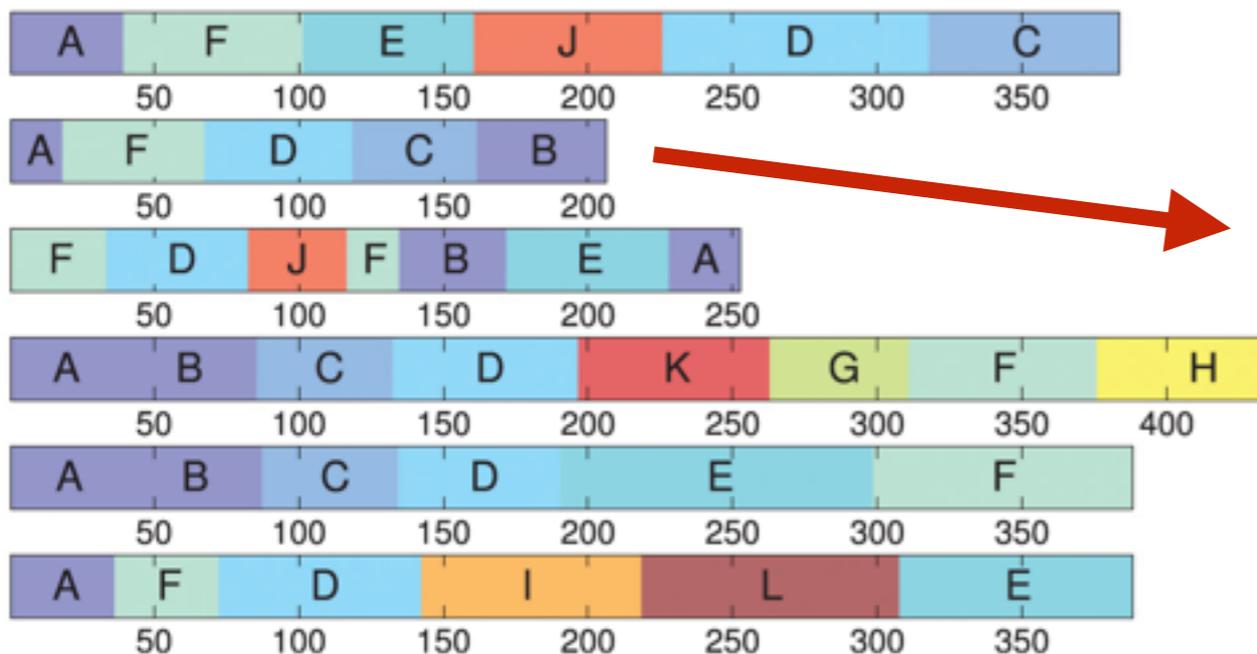
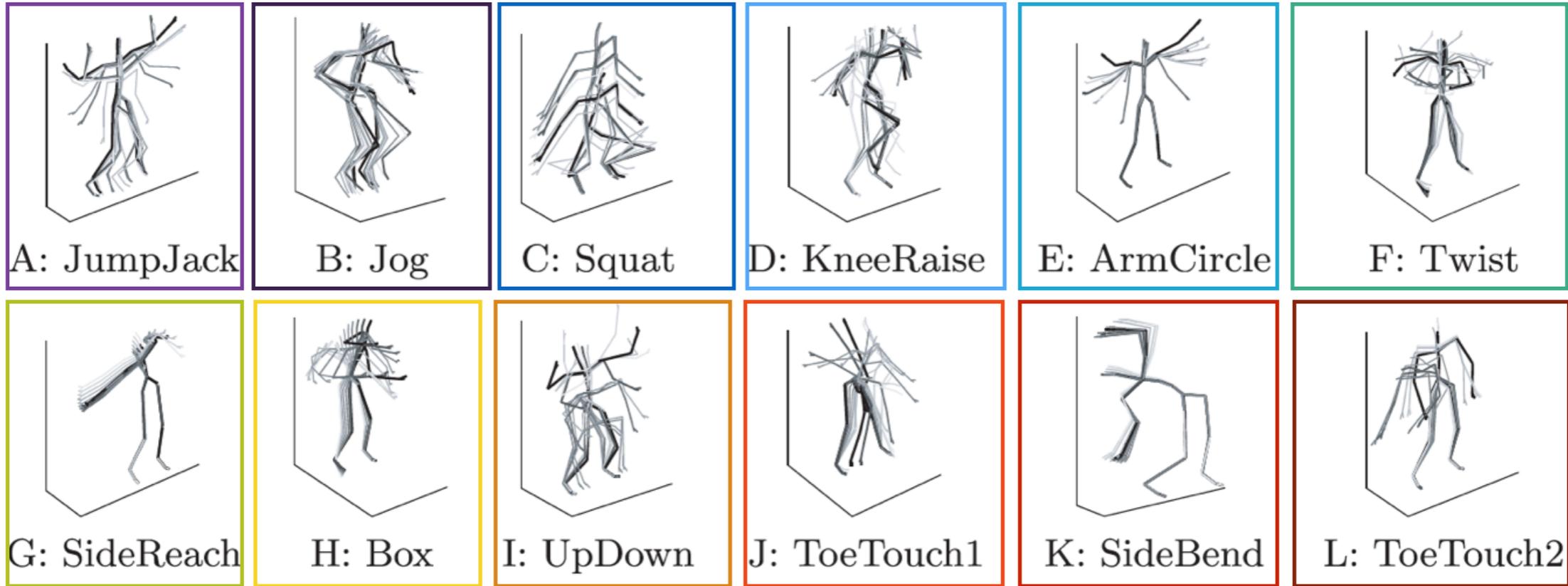
Hierarchical ddCRP



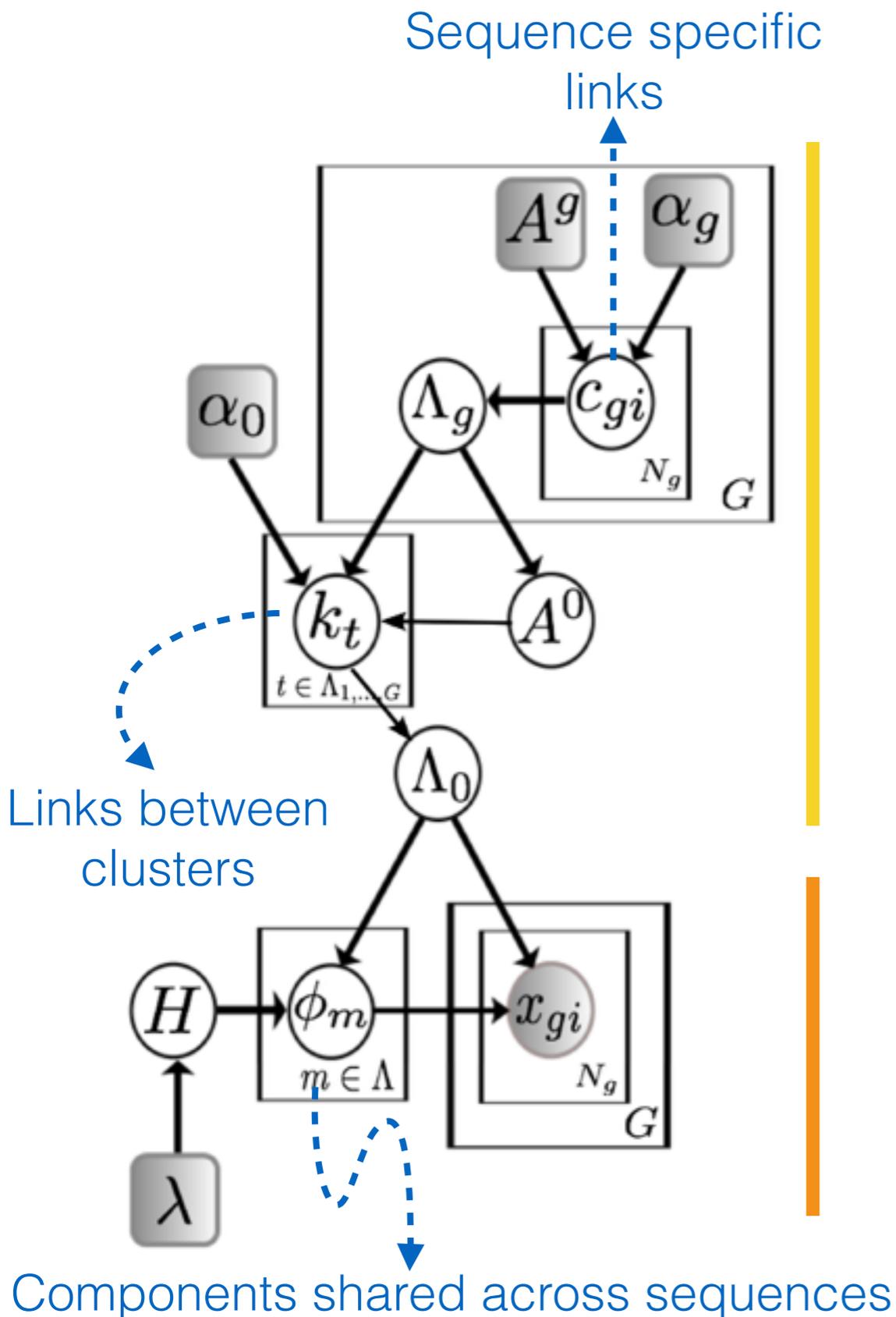
Hierarchical ddCRP



Activity Recognition



Hierarchical Auto Regressive Mixtures



- Sequence specific ddCRP models:

$$p(c_{gi} = gj \mid \alpha_g, A^g) \propto \begin{cases} \exp(-\frac{(i-j)}{N_g^\gamma}) & i > j, \\ 0 & i < j, \\ 1 & i = j. \end{cases}$$

- Global CRP across sequences:

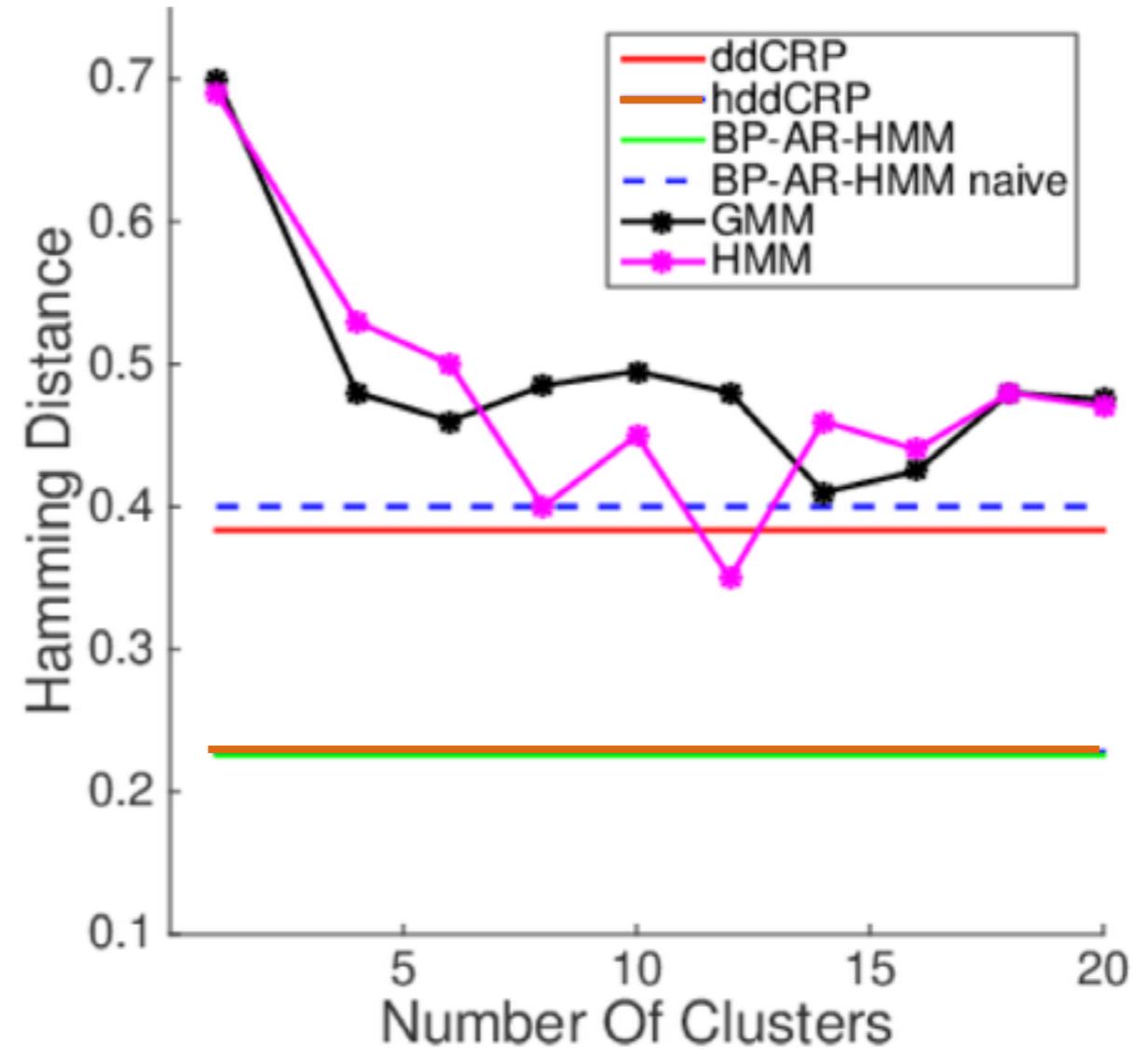
$$\Lambda_0 \sim \text{CRP}(T(\mathbf{c}), \alpha_0)$$

- Autoregressive likelihoods:

$$x_{gt} = B_m x_{gt-1} + \epsilon_m$$

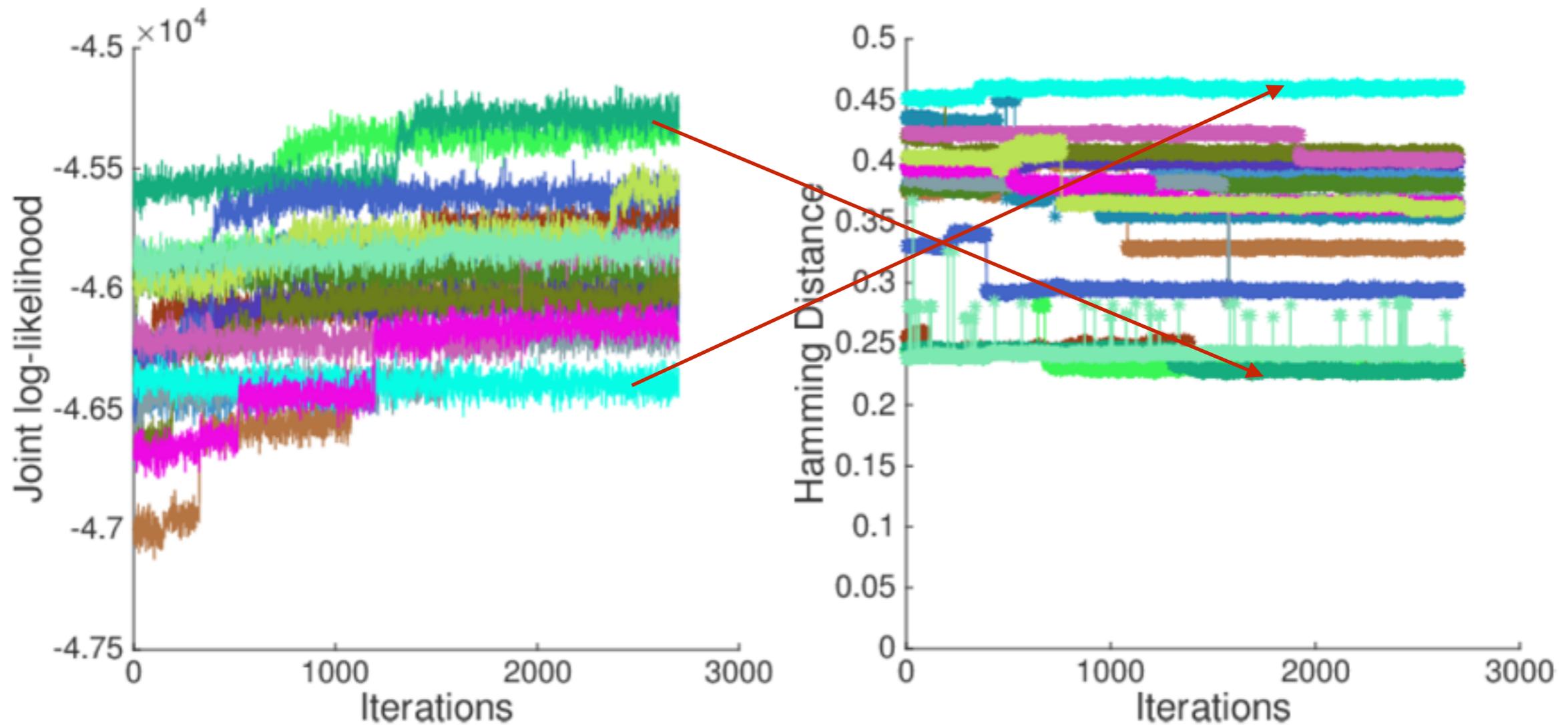
$$B_m, \epsilon_m \sim H(\lambda)$$

Discovered Activities



Examples of activities discovered by hddCRP

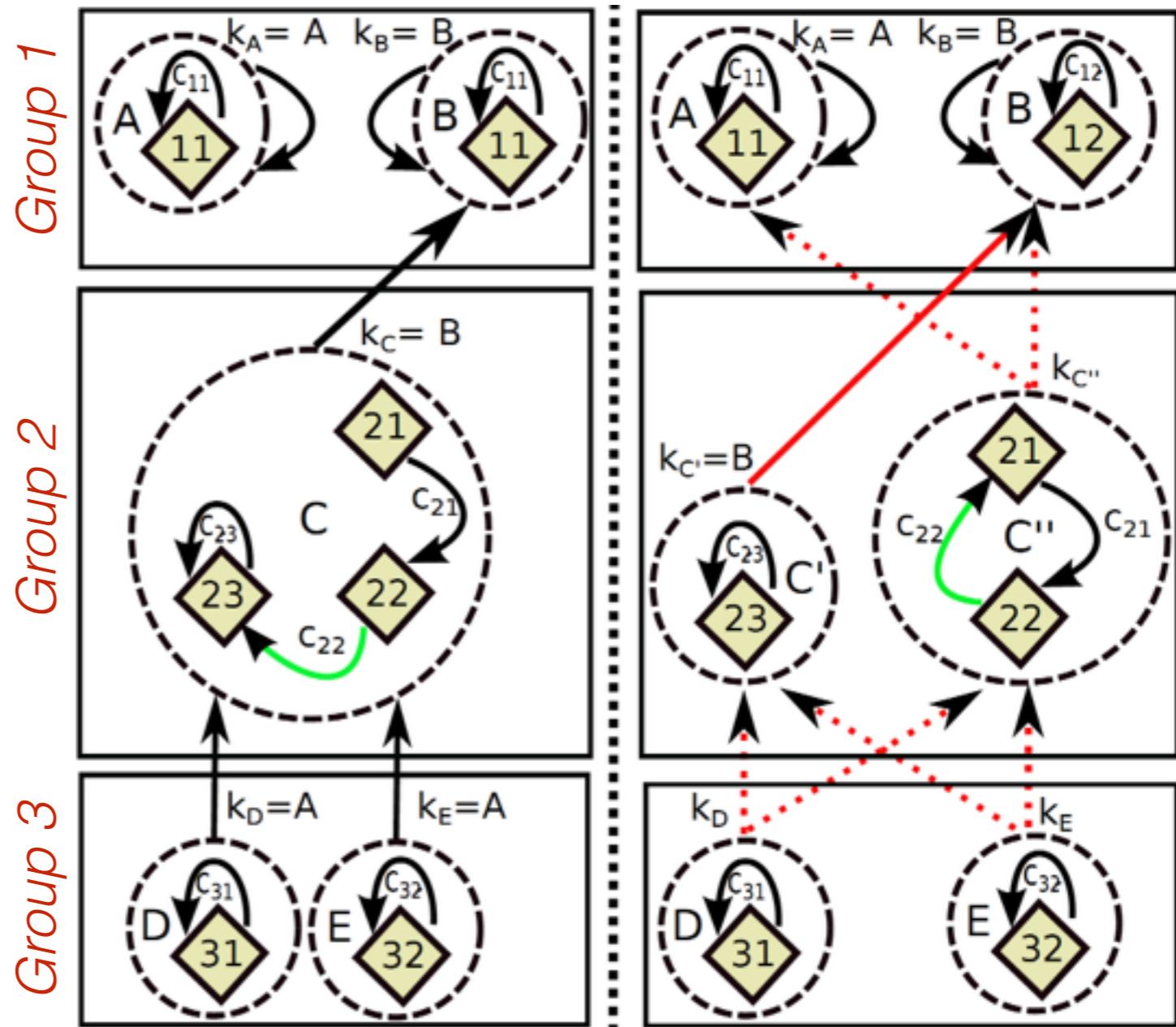
External Model Validation



15 independent MCMC chains

Inference

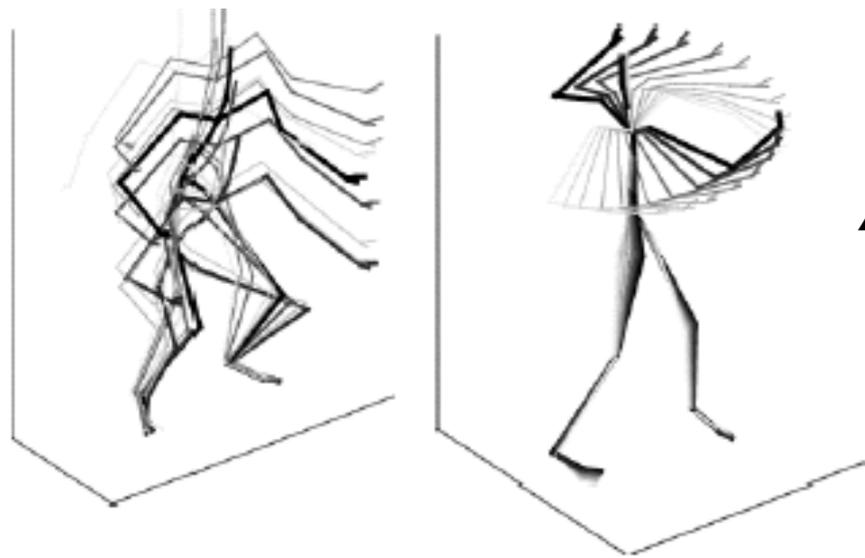
- More involved.
- No simple Gibbs sampler, need to resort to Metropolis Hastings.
- Nonetheless “efficient” MH samplers can be crafted.



Summary



*Articulated object segmentation
through ddCRP mixtures*



*Activity discovery via hierarchical
distance dependent models*

Talk Outline

- Distance dependent partitions
 - Parts from articulated 3D objects
- Hierarchical distance dependent partitions
 - Activity discovery from MoCap data
- Learning distance dependent models
 - Image and video segmentation

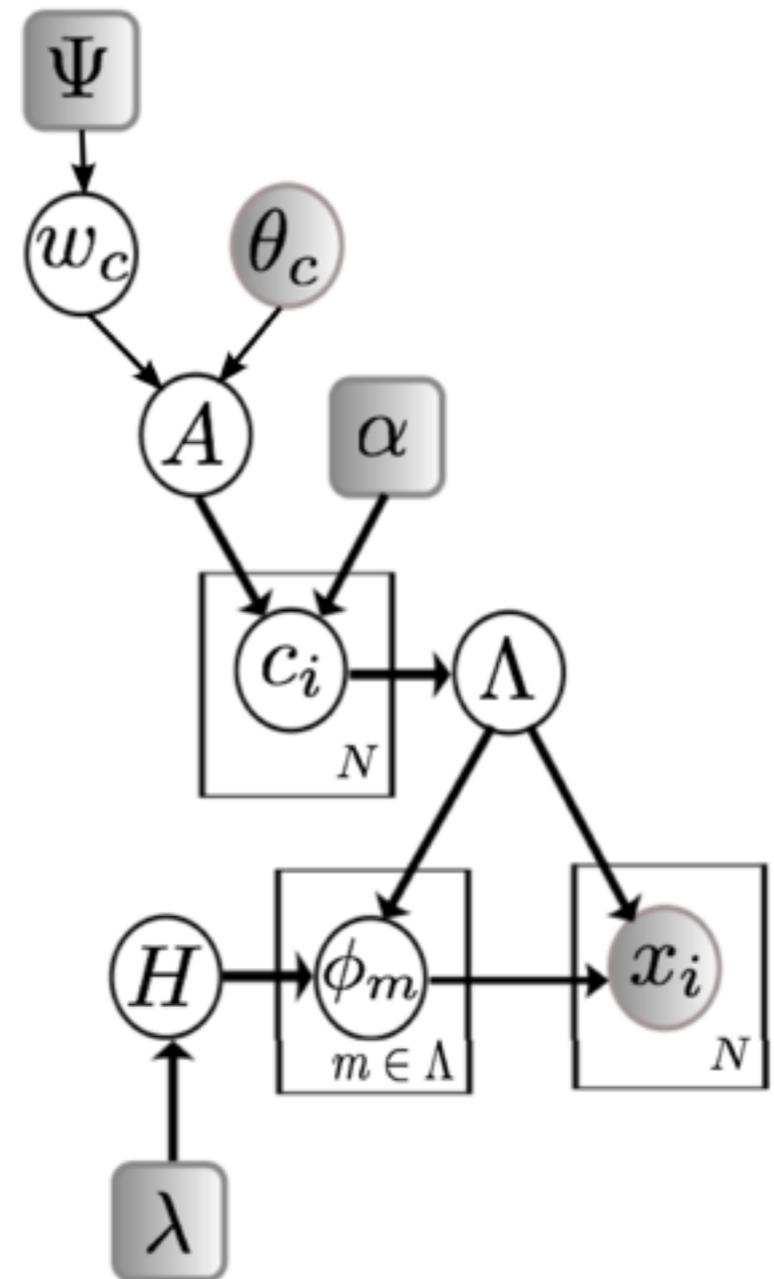
Feature Augmented Models

$$p(c_i = j | A) \propto A_{ij}$$

$$A_{ij} = f(w_c^T \theta_{ij}^c)$$

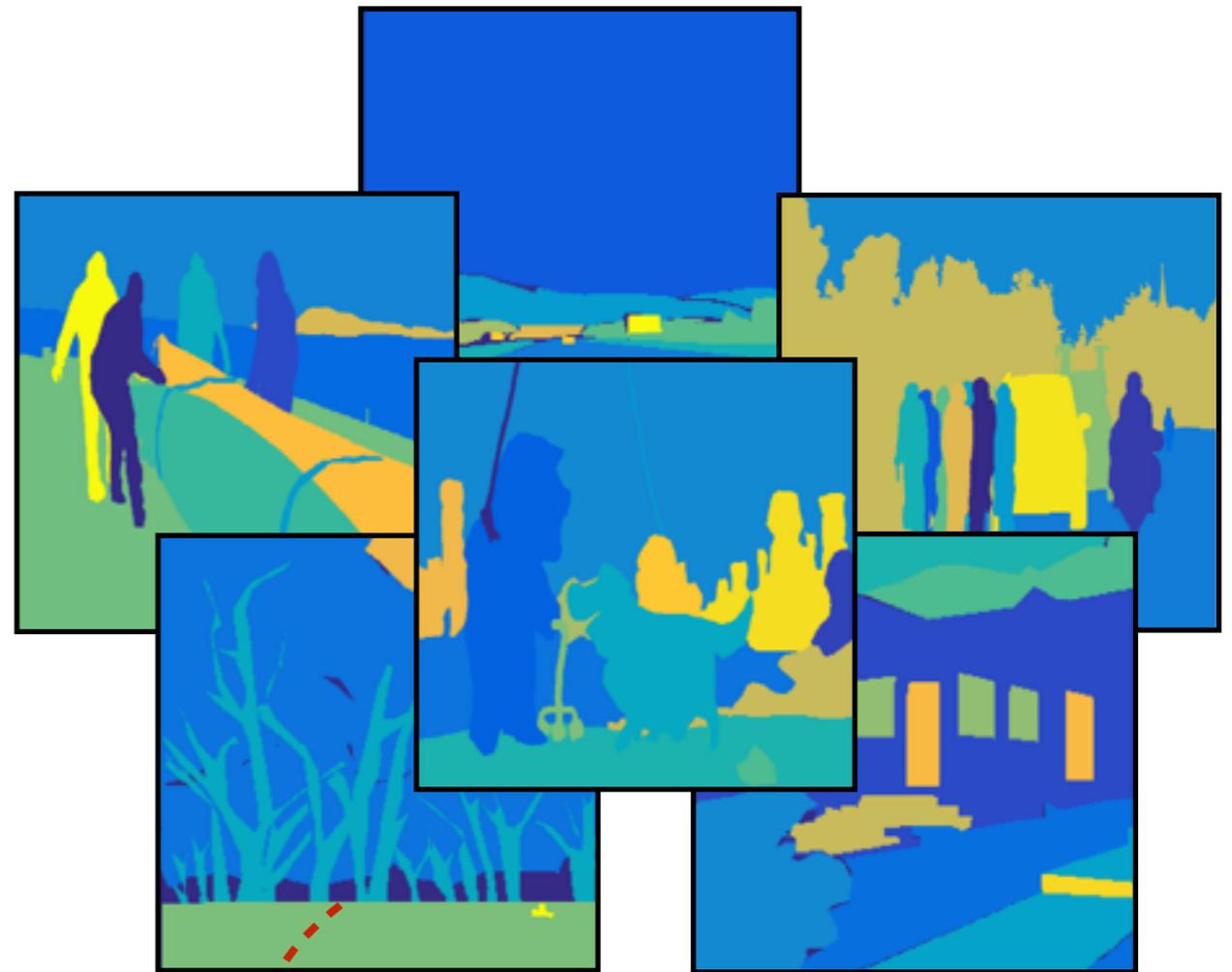
Features encoding
similarity

Latent variables
governing contribution
of features



Learning From Partitions

- Moderate sized databases of partitions available for image and video collections.
- *Uncertainty* in labeled partitions
- Partitions are observed, but *links are not*.



$$y_i = \mathbb{N}^{N_k \times 1} = [1, 1, 2, 4, \dots, 3, 3]$$

$$Y = \{y_1 \dots y_D\}$$

Approximate Bayesian Computation

- Noisy partitions - human interpretations vary
 - Appropriate noise model? Unclear, *ABC* *instead*
- Likelihood free inference:
 - Match “*interesting*” model statistic with observed data statistic

Auxiliary Training Model

$$p(\mathbf{c}, w_c, Y) \propto p(w_c) \prod_{d=1}^D p(c_d | w_c) \mathbf{1}(z(\mathbf{c}_d), y_d)$$

$$\mathbf{1}(y_a, y_b) = \begin{cases} 1 & \text{if } \Delta(y_a, y_b) < \epsilon, \\ 0 & \text{otherwise.} \end{cases}$$

*Probability restricted to partitions **close** to training data.*

Loss Aware Model

- Notion of closeness captured through a task specific loss function:

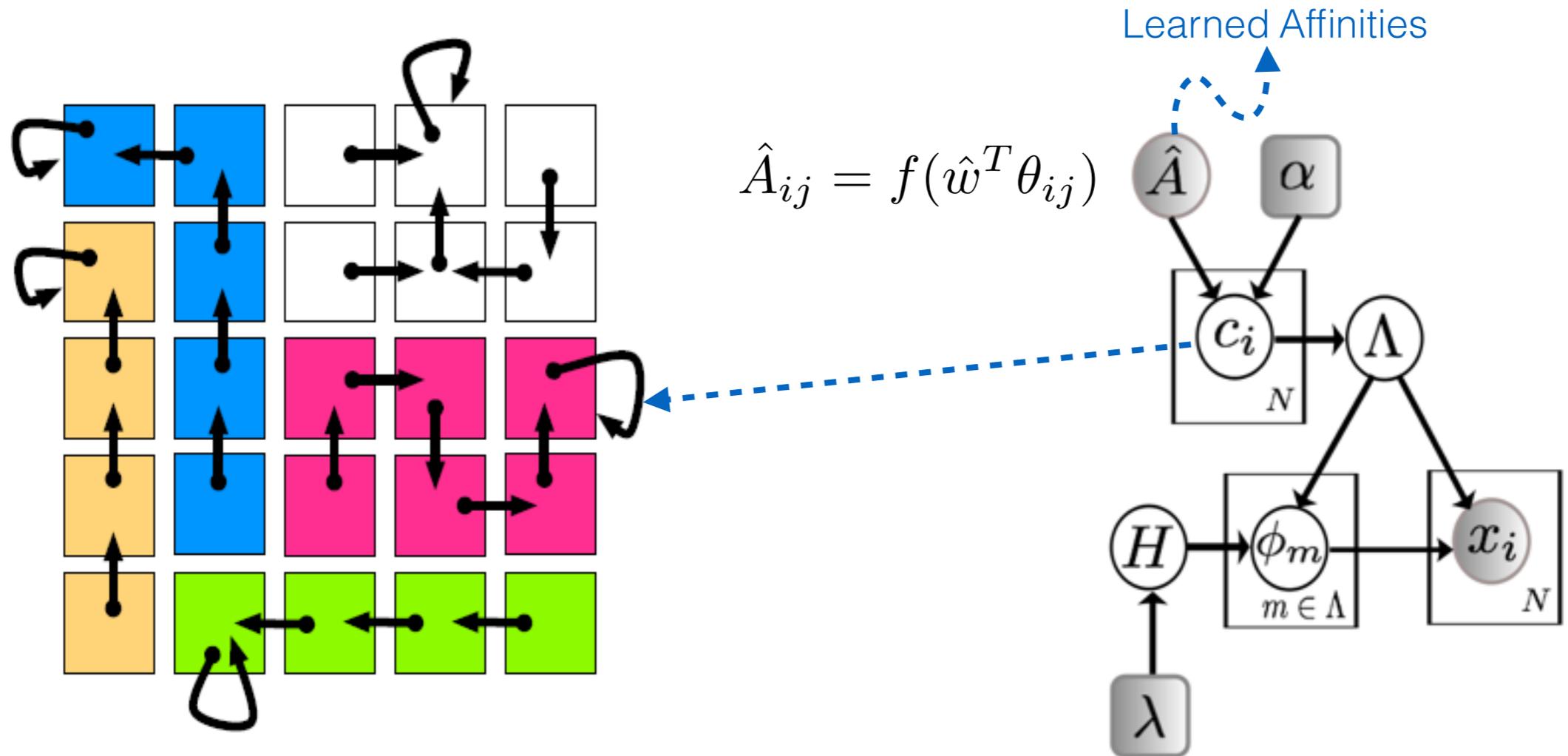
$$\Delta(y_a, y_b) = 1 - \text{RI}(y_a, y_b)$$

- *Marginalize* over the exponentially large space of *latent links* using MCMC
- Efficient ABC variant for sampling from the *auxiliary training model*

Talk Outline

- Distance dependent partitions
 - Parts from articulated 3D objects
- Hierarchical distance dependent partitions
 - Activity discovery from MoCap data
- Learning distance dependent models
 - Image and video segmentation

Image Segmentation



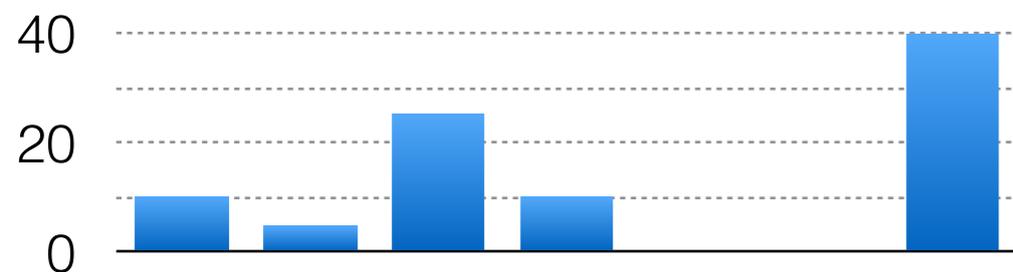
Generative features:

$$\theta_{ij} = \{\text{row}_i - \text{row}_j, \\ \text{col}_i - \text{col}_j\}$$

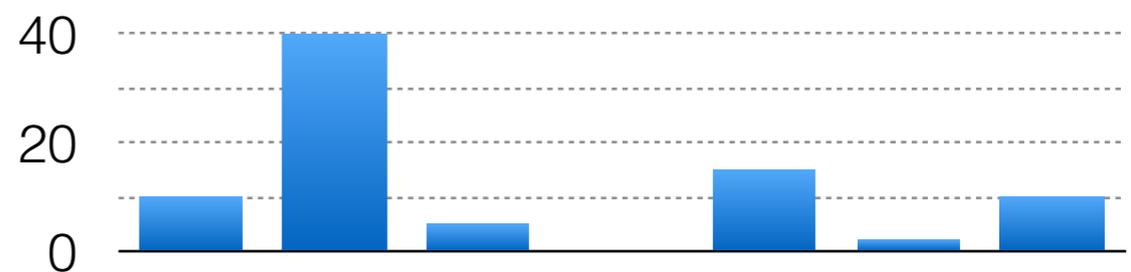
Conditional features:

$$\theta_{ij} = \{\text{row}_i - \text{row}_j, \\ \text{col}_i - \text{col}_j, \text{edge}_{ij}\}$$

Image Representation



$$x_i^{color} \sim \text{Mult}(\phi_m^{color})$$



$$x_i^{texture} \sim \text{Mult}(\phi_m^{texture})$$

Each super-pixel is described through histograms (~120 bin) of color and texture

Eight Natural Scene Category Dataset (LabelMe)



400 train and 800 test images

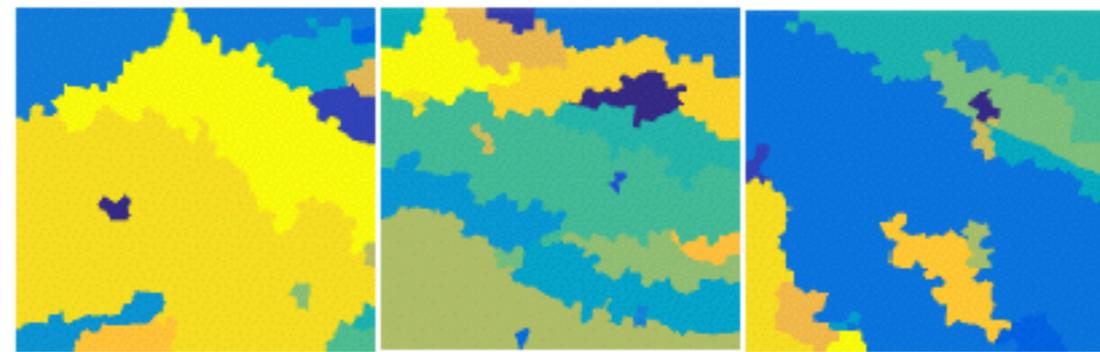
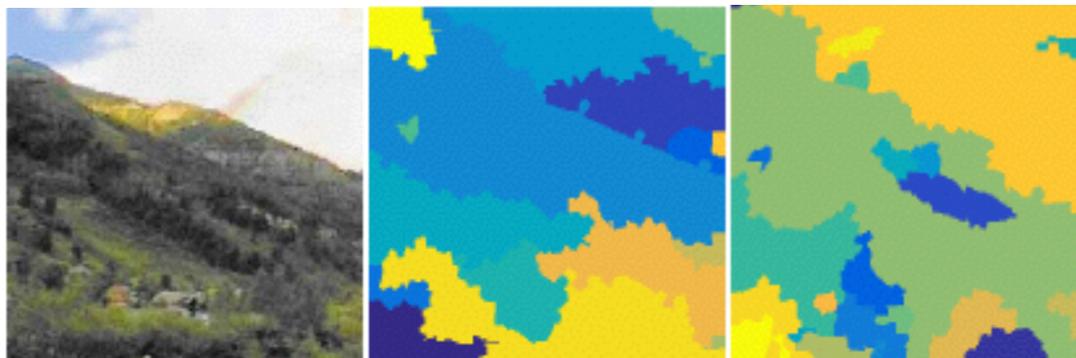
Oliva and Torralba, 2001

Samples from learned models

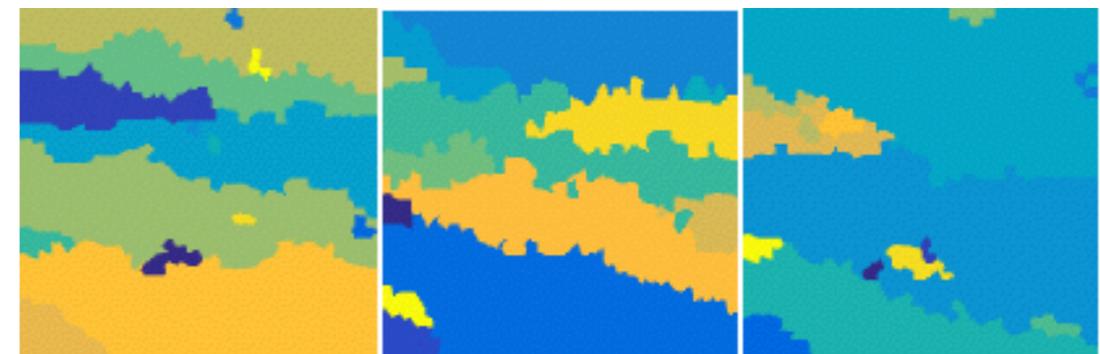
Conditional

Generative

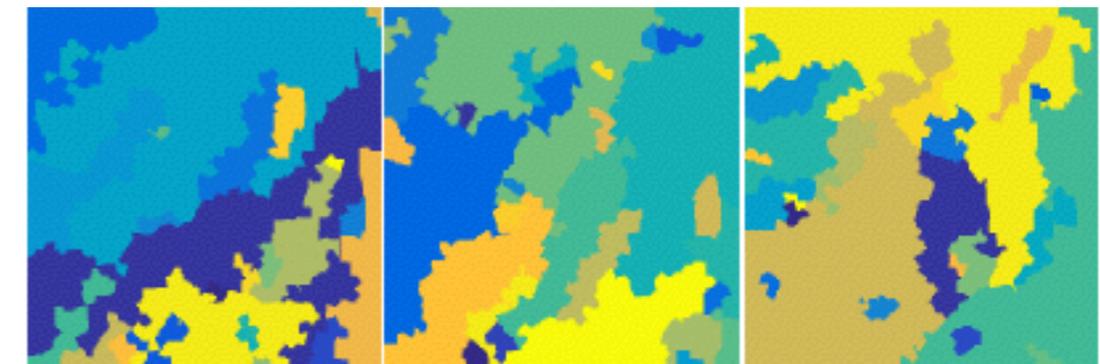
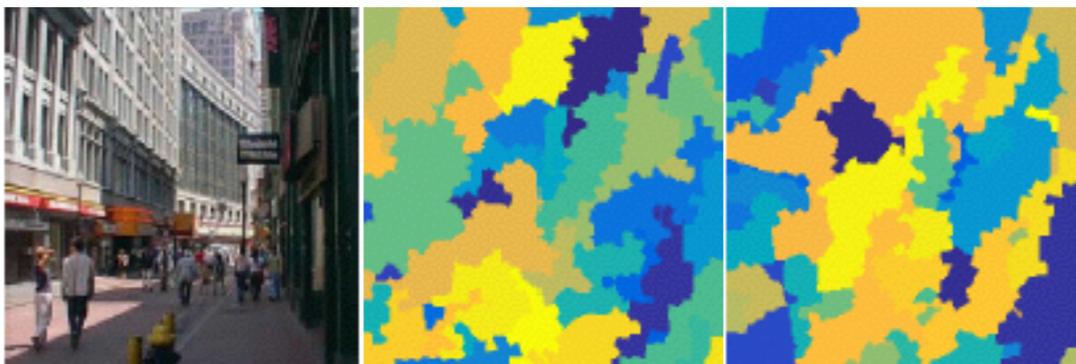
Mountain



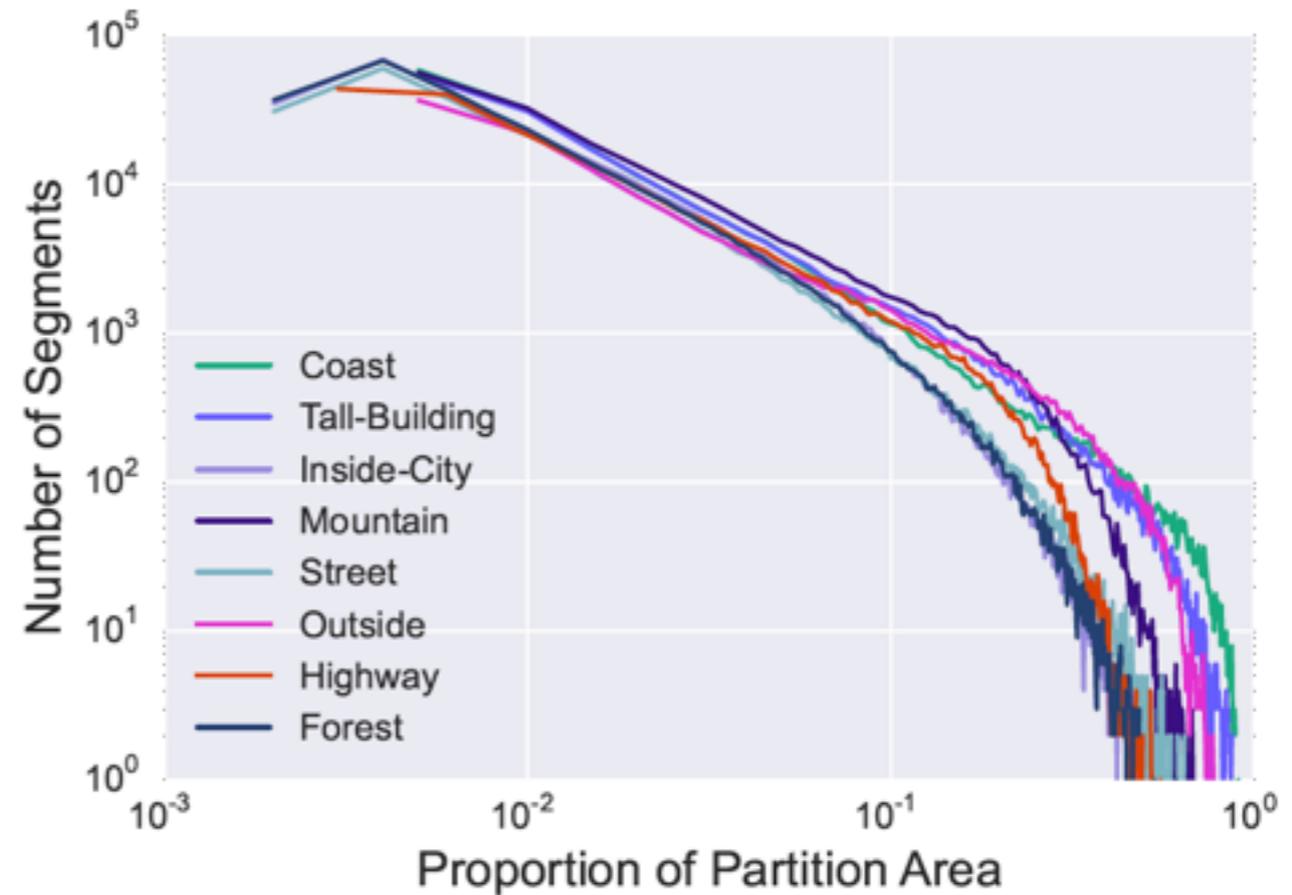
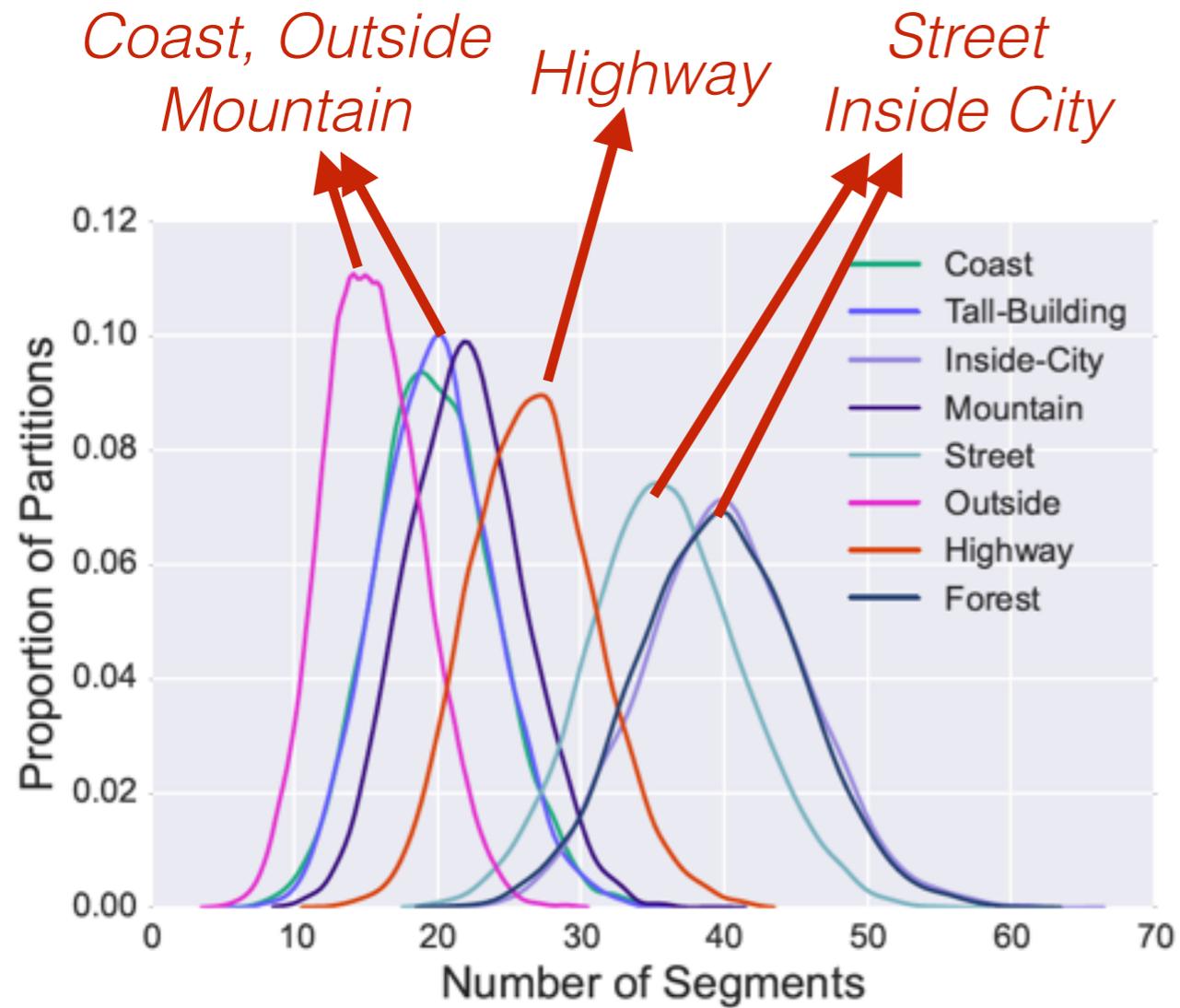
Coast



Street

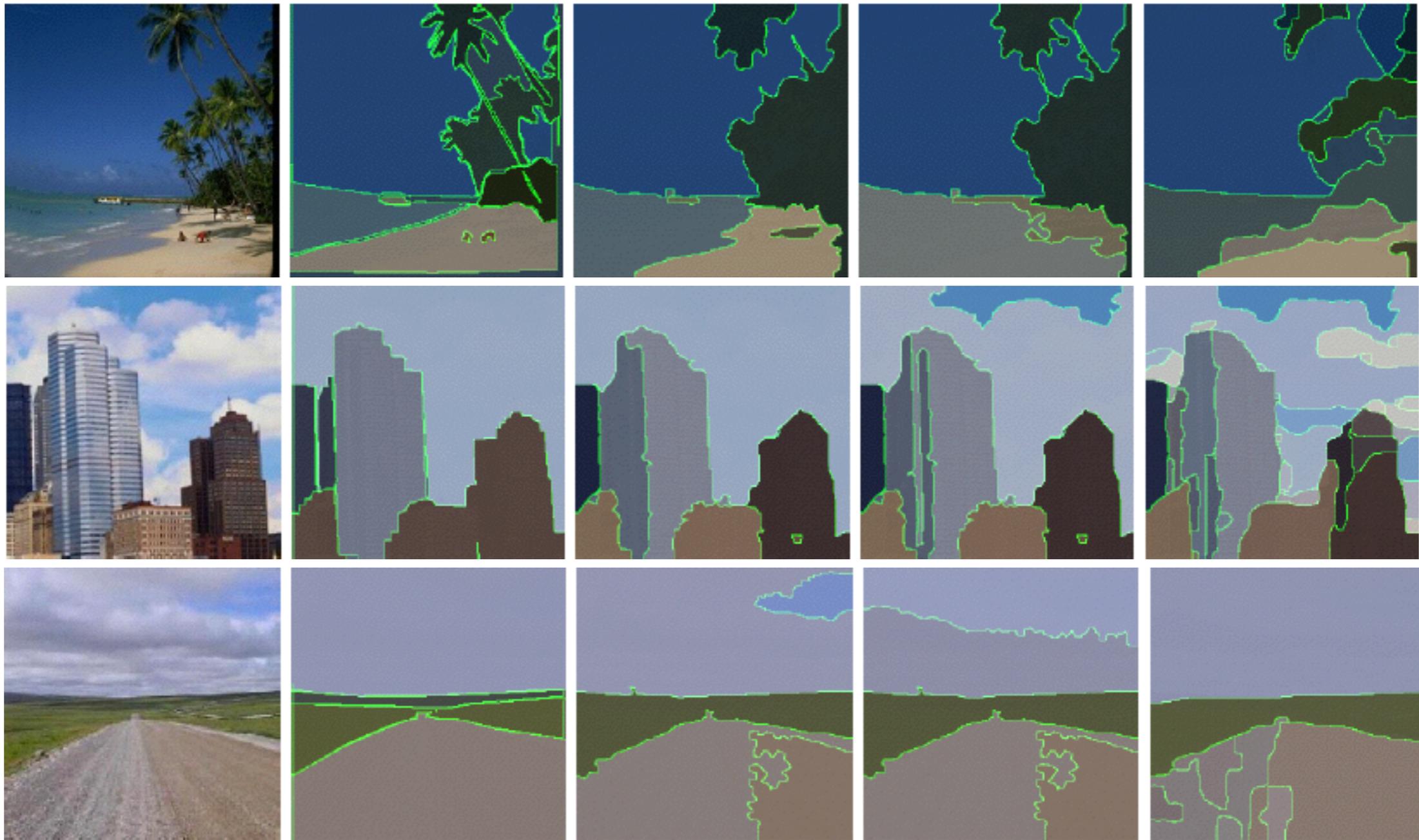


Monte Carlo Statistics



Statistics from 10,000 partitions sampled from generative affinities

Qualitative Results



GT

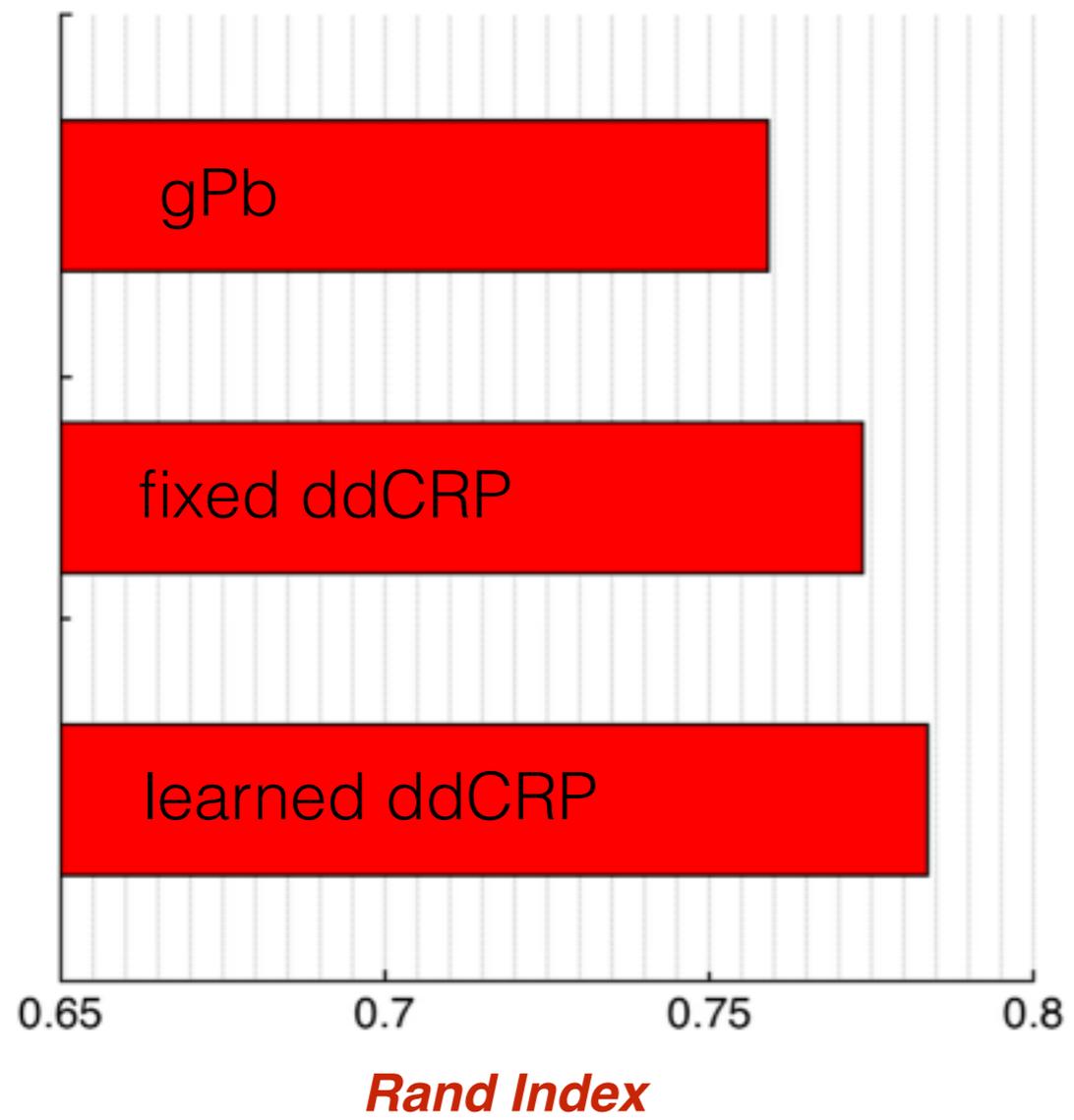
learned
ddCRP

fixed
ddCRP

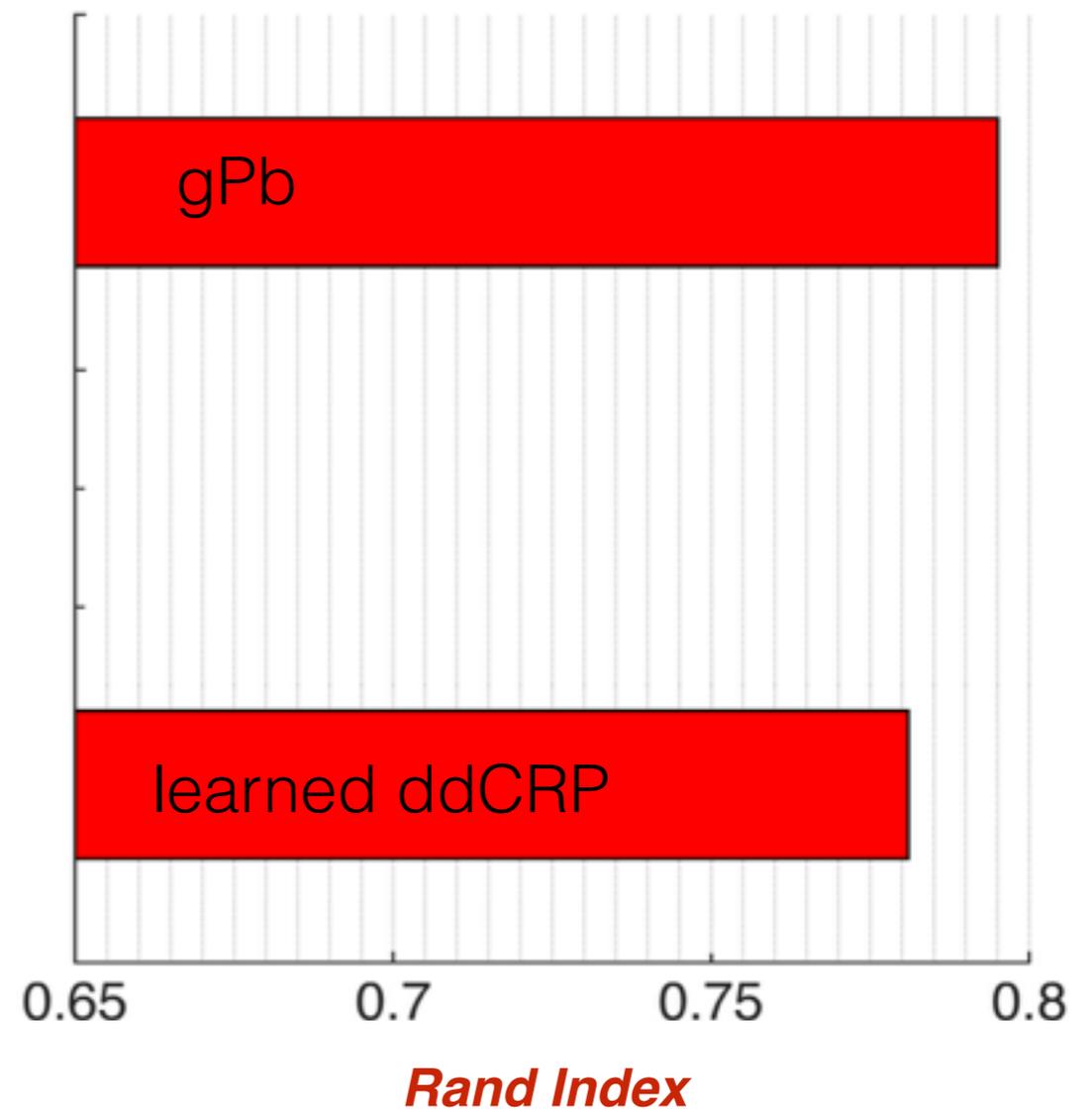
gPb

Quantitative results

LabelMe



BSDS300

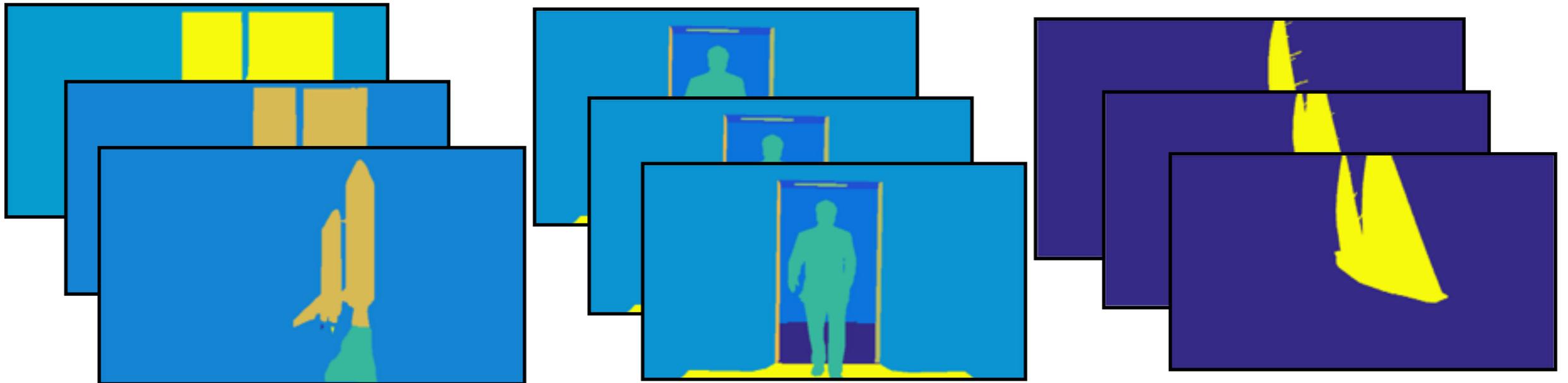


Learning in hierarchical models

- Auxiliary model for training now needs to account for links between clusters

$$p(\mathbf{c}, \mathbf{k}, w, Y) \propto p(w) \prod_{d=1}^D p(c_d | w_c) p(k_d | c_d, w_k) \mathbf{1}(z(\mathbf{c}_d, k_d), y_d)$$

$$w = \{w_c, w_k\}$$



VSB 100 - 40 training videos

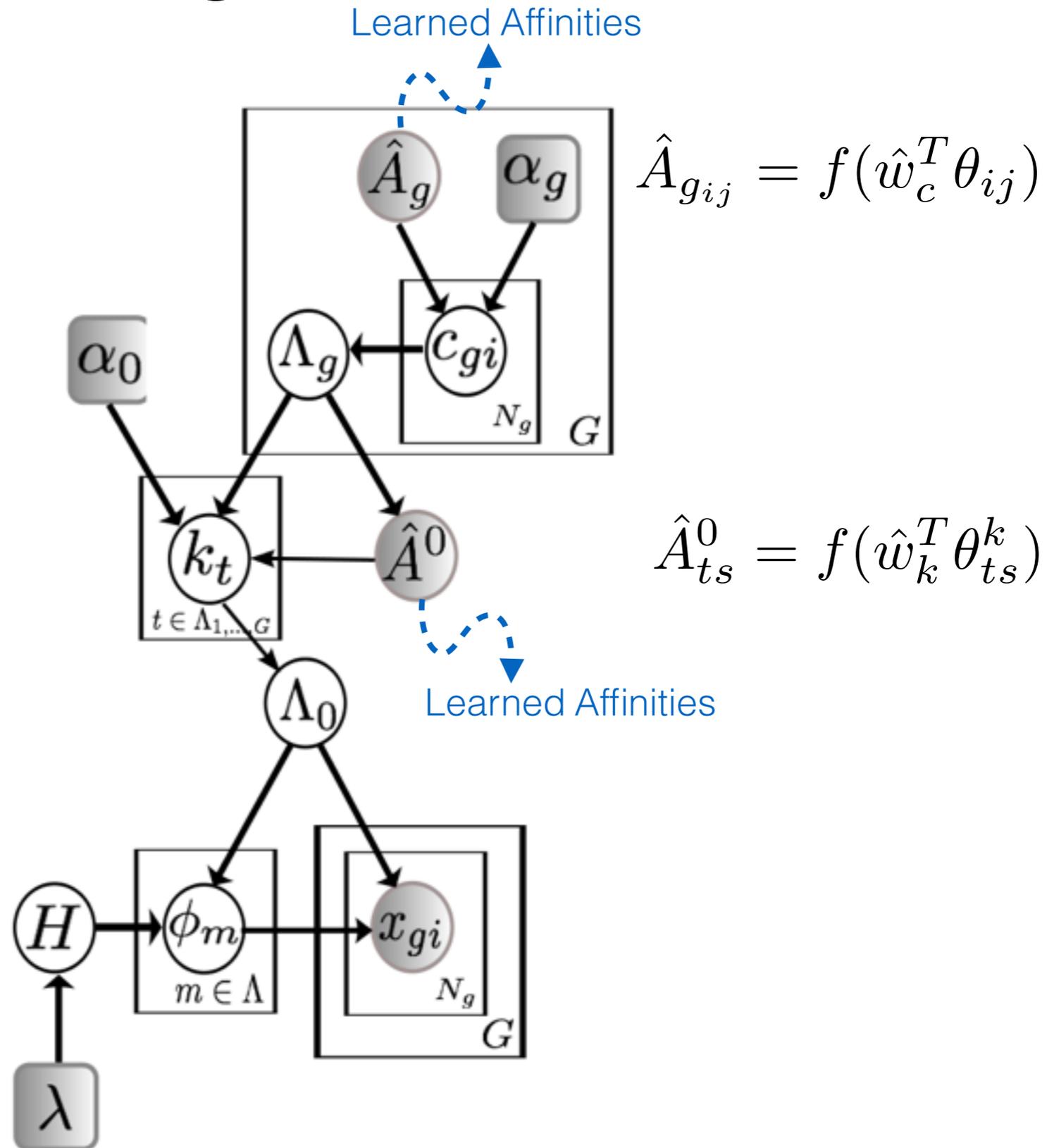
Video Segmentation

Features between superpixels:

$$\theta_{ij} = \{ \text{row}_i - \text{row}_j, \\ \text{col}_i - \text{col}_j, \\ \text{edge}_{ij} \}$$

Features encoding similarity between segments:

$$\theta_{ts}^k = \{ \psi(\text{size}_{ts}, \\ \text{shape}_{ts}, \\ \text{locations}_{ts}) \}$$



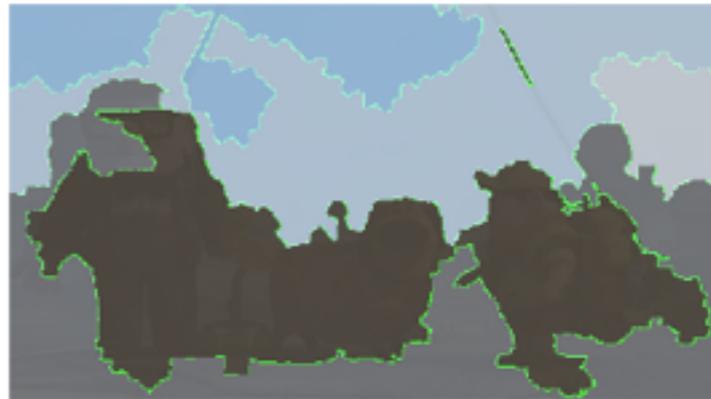
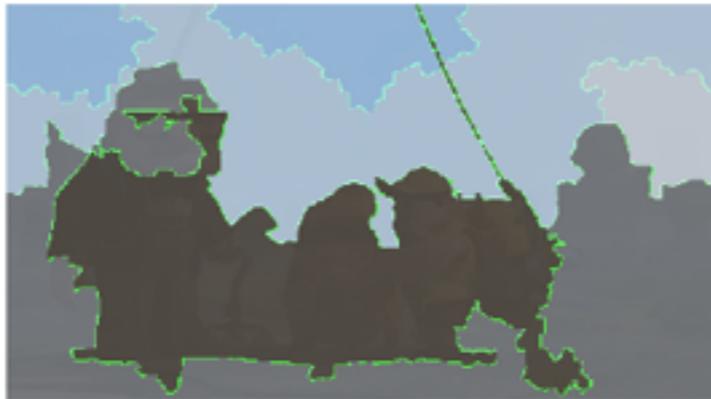
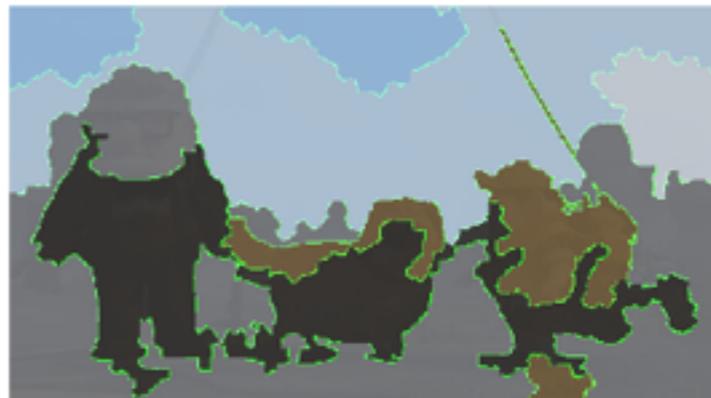
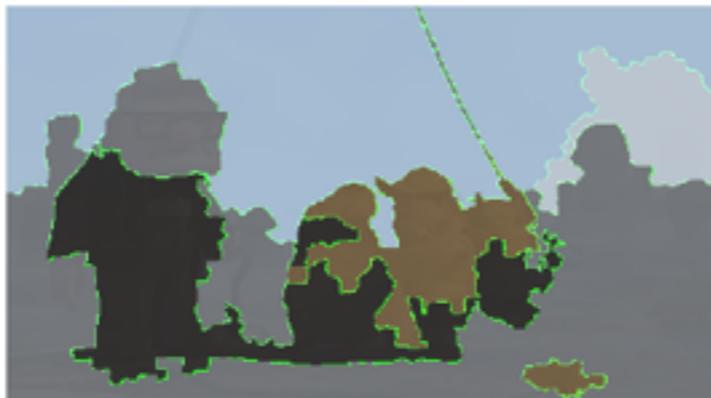
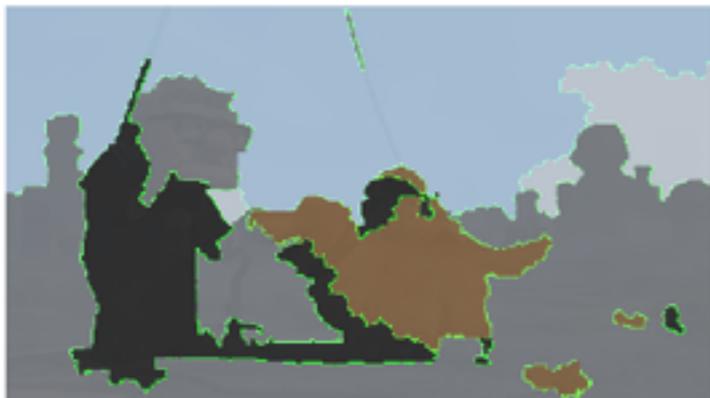
First Frame

Last Frame

GT

learned

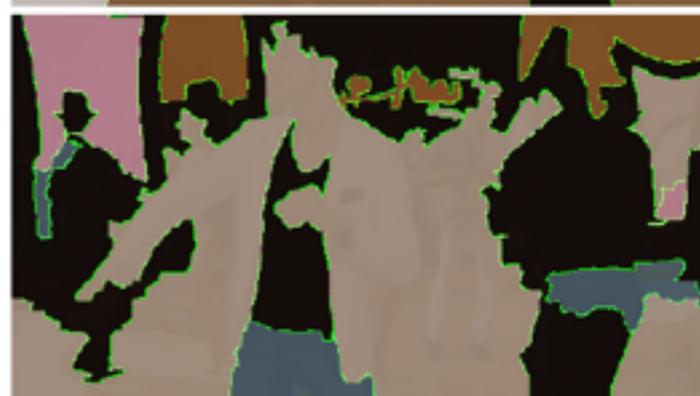
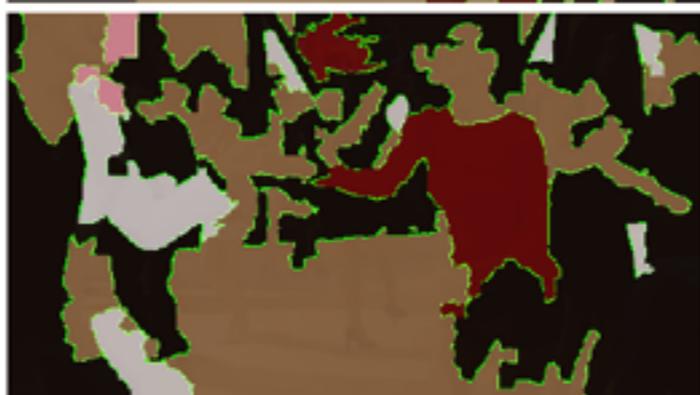
fixed



First Frame

Last Frame

GT
learned
fixed



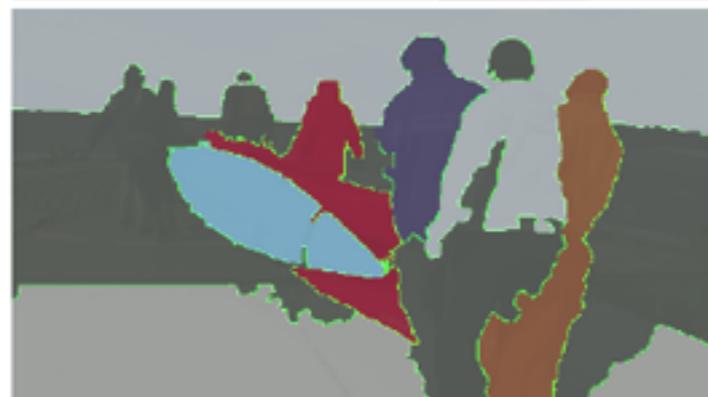
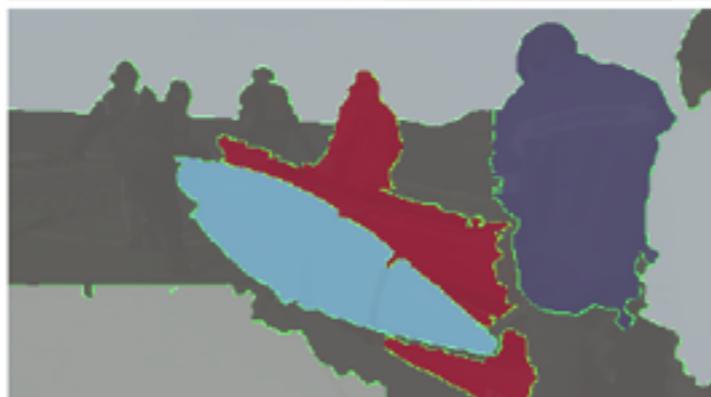
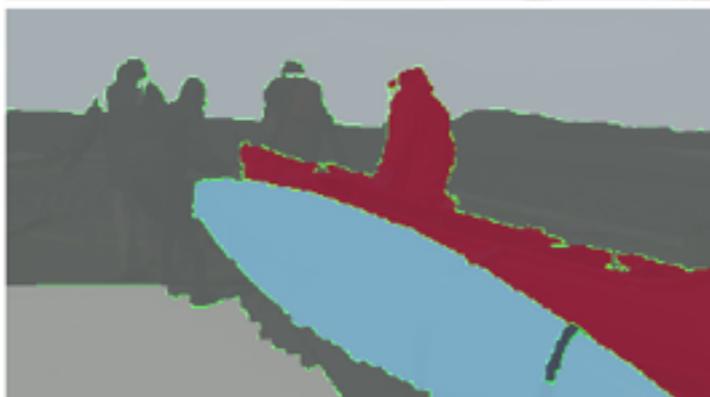
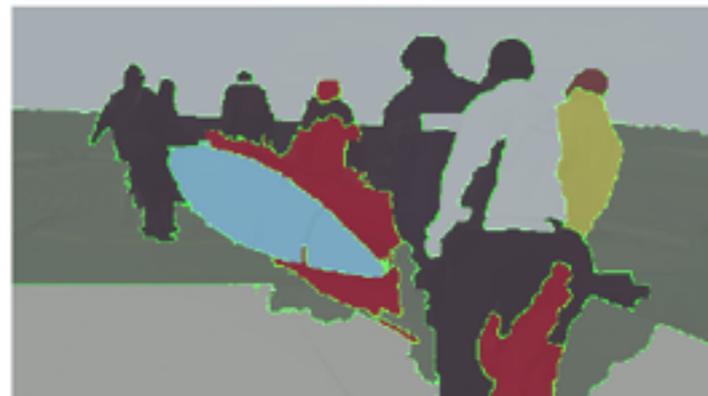
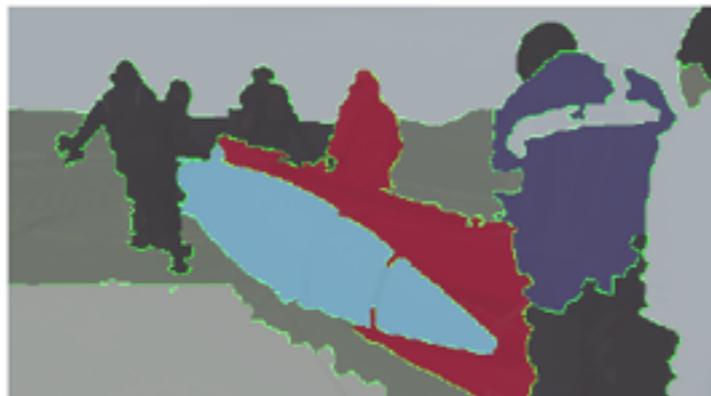
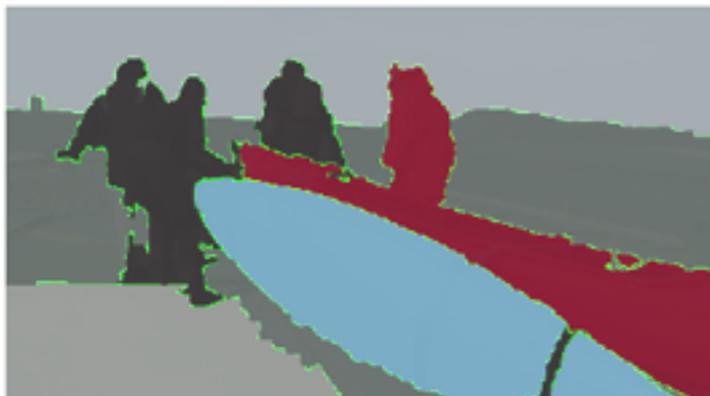
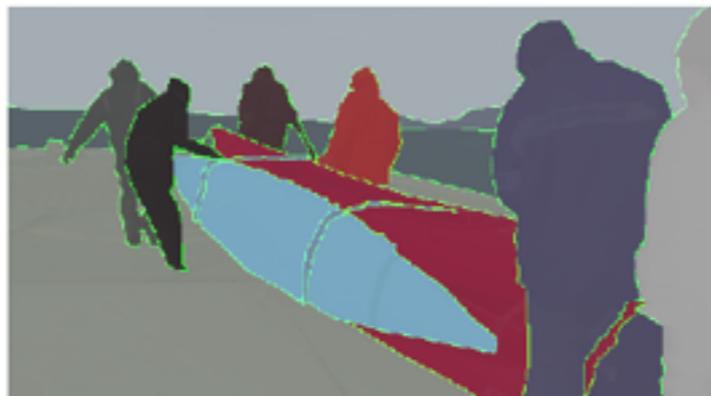
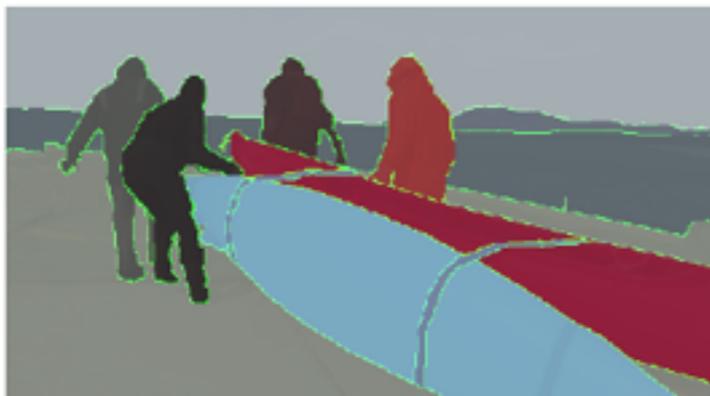
First Frame

Last Frame

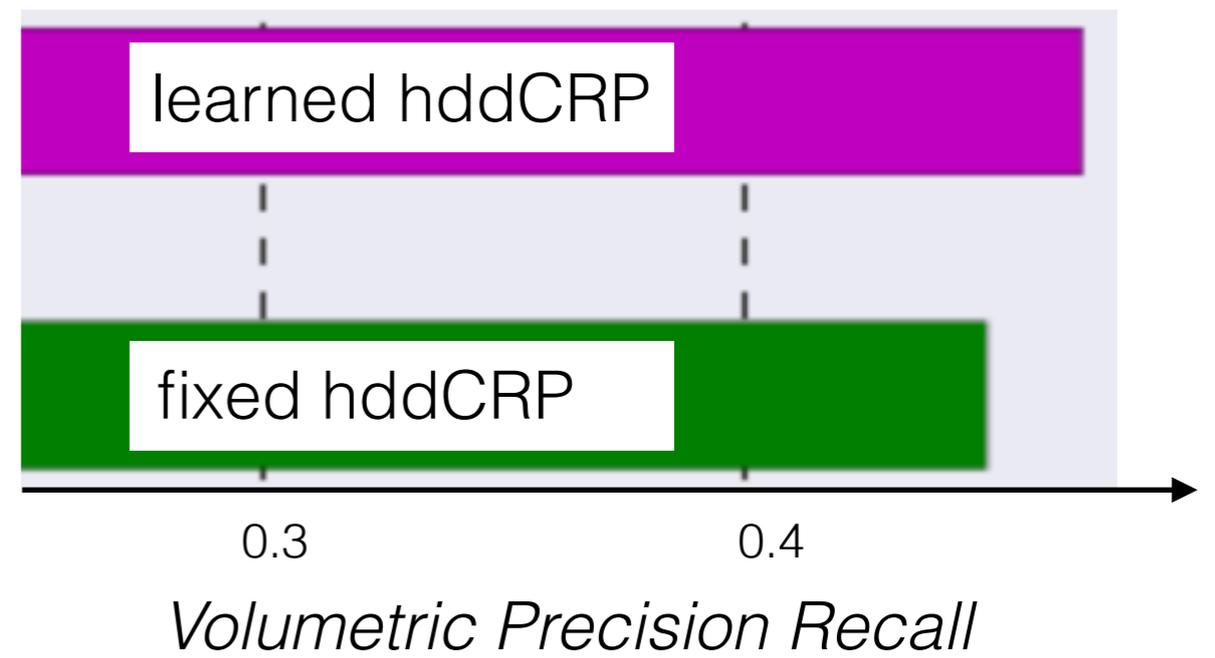
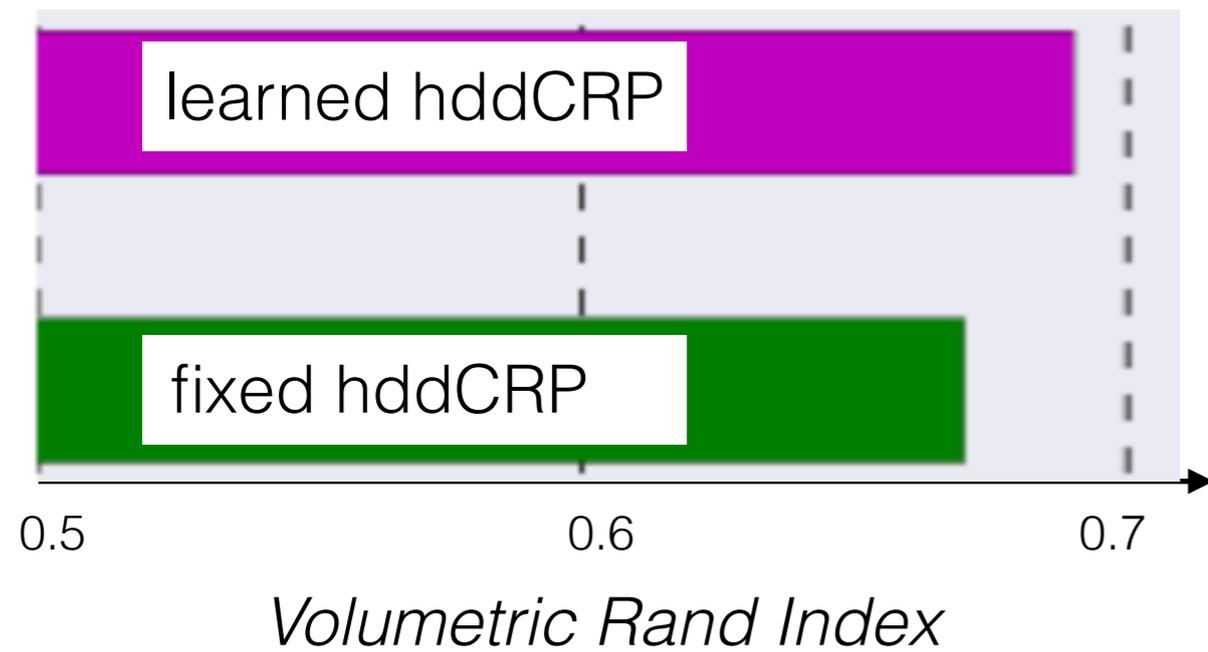
GT

learned

fixed



Learning benefits hddCRP



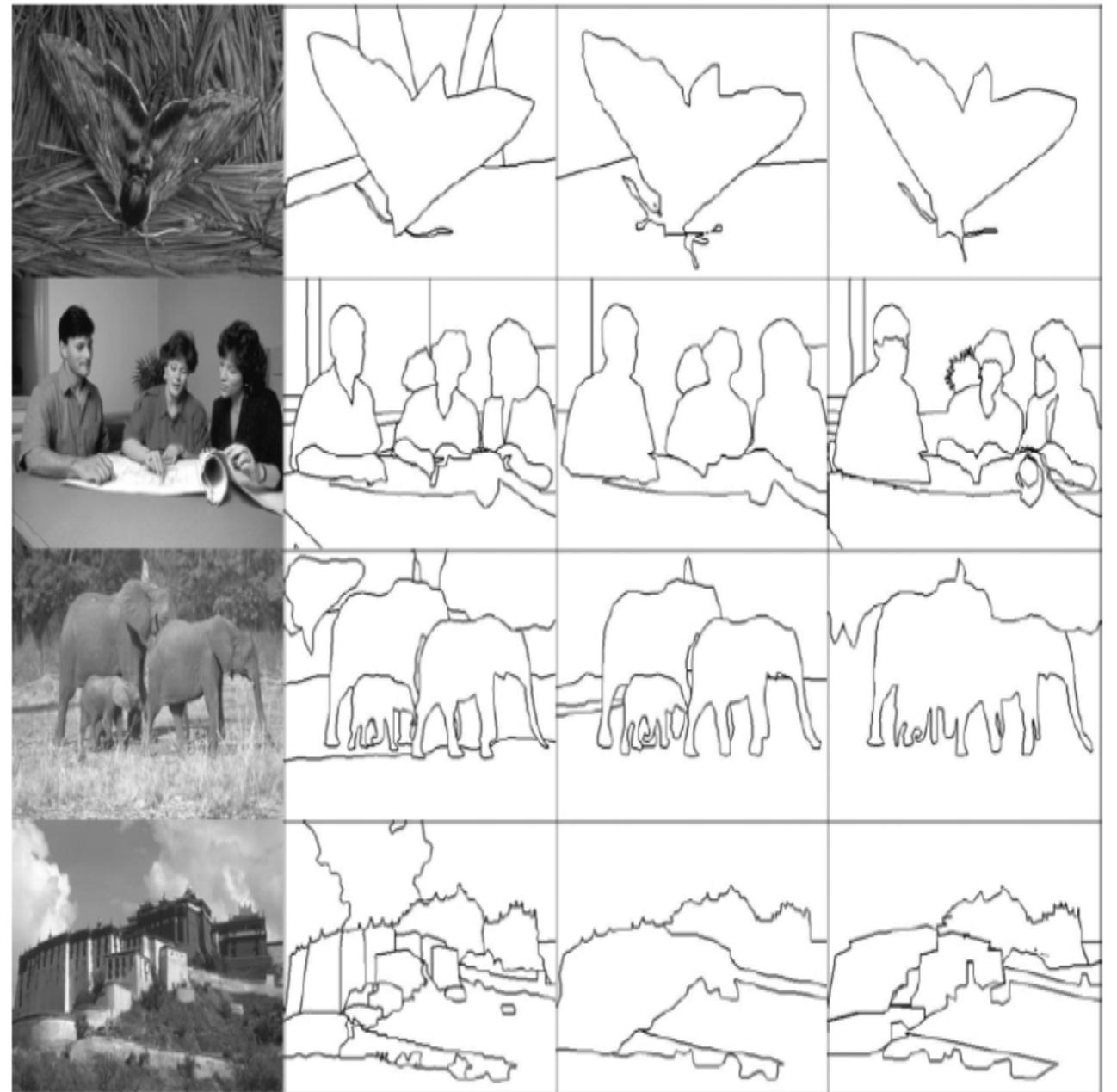
Thank You



Questions?

Statistics of Human Segments

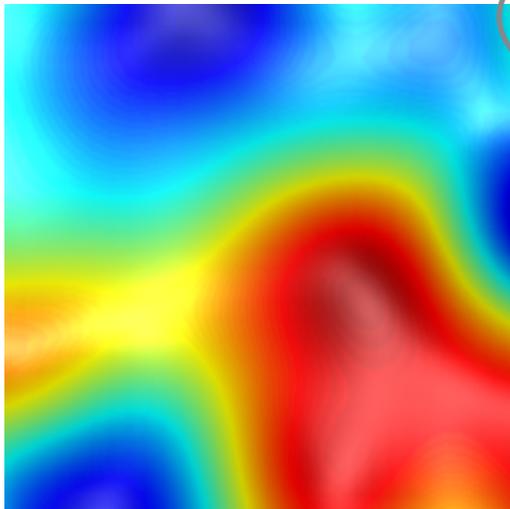
- Human segment sizes follow power law behavior.



Spatial Coupling through Layers

Smooth Layers

$$u_1 \sim \text{GP}(0, K)$$



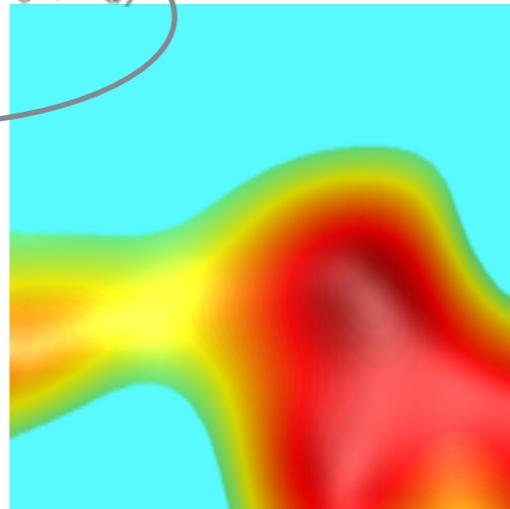
$$w_1 \sim \text{Beta}(1 - \alpha_a, \alpha_b + \alpha_a)$$

$$\delta_1 = \Phi^{-1}(w_1)$$

$$u_1 < \delta_1$$



Thresholded layer support



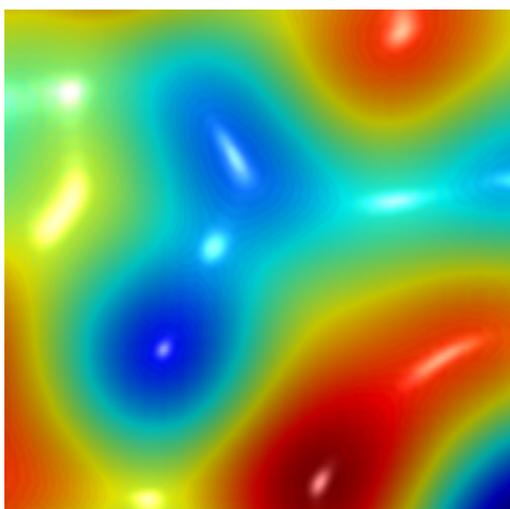
Occlude

Image Partition



$$z_n = \min\{k \mid u_{kn} < \delta_k\}$$

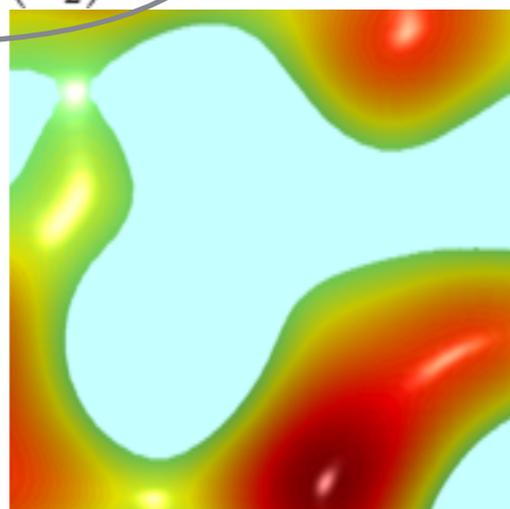
$$u_2 \sim \text{GP}(0, K)$$



$$w_2 \sim \text{Beta}(1 - \alpha_a, \alpha_b + 2\alpha_a)$$

$$\delta_2 = \Phi^{-1}(w_2)$$

$$u_2 < \delta_2$$



Sudderth & Jordan, 2008

Ghosh & Sudderth, 2012

Video Segmentation

- Features between superpixels — same as image segmentation.
- Features between segments — Shapes, sizes and positions.

$$\theta_{ts}^k = [\vartheta_{ts}, \varphi_{ts}, \frac{|\zeta_t - \zeta_s|}{S}]^T,$$

$$\vartheta_{ts} = \mathbf{1}_{[t,s|t \in g, s \in g]} \left[\frac{r_t - r_s}{R}, \frac{y_t - y_s}{Y} \right]^T,$$

$$\varphi_{ts} = \mathbf{1}_{[t,s|t \in g+1, s \in g]} \left[\frac{|r_t - r_s|}{R}, \frac{|y_t - y_s|}{Y}, 1 - \frac{t \cap s}{t \cup s} \right]^T$$

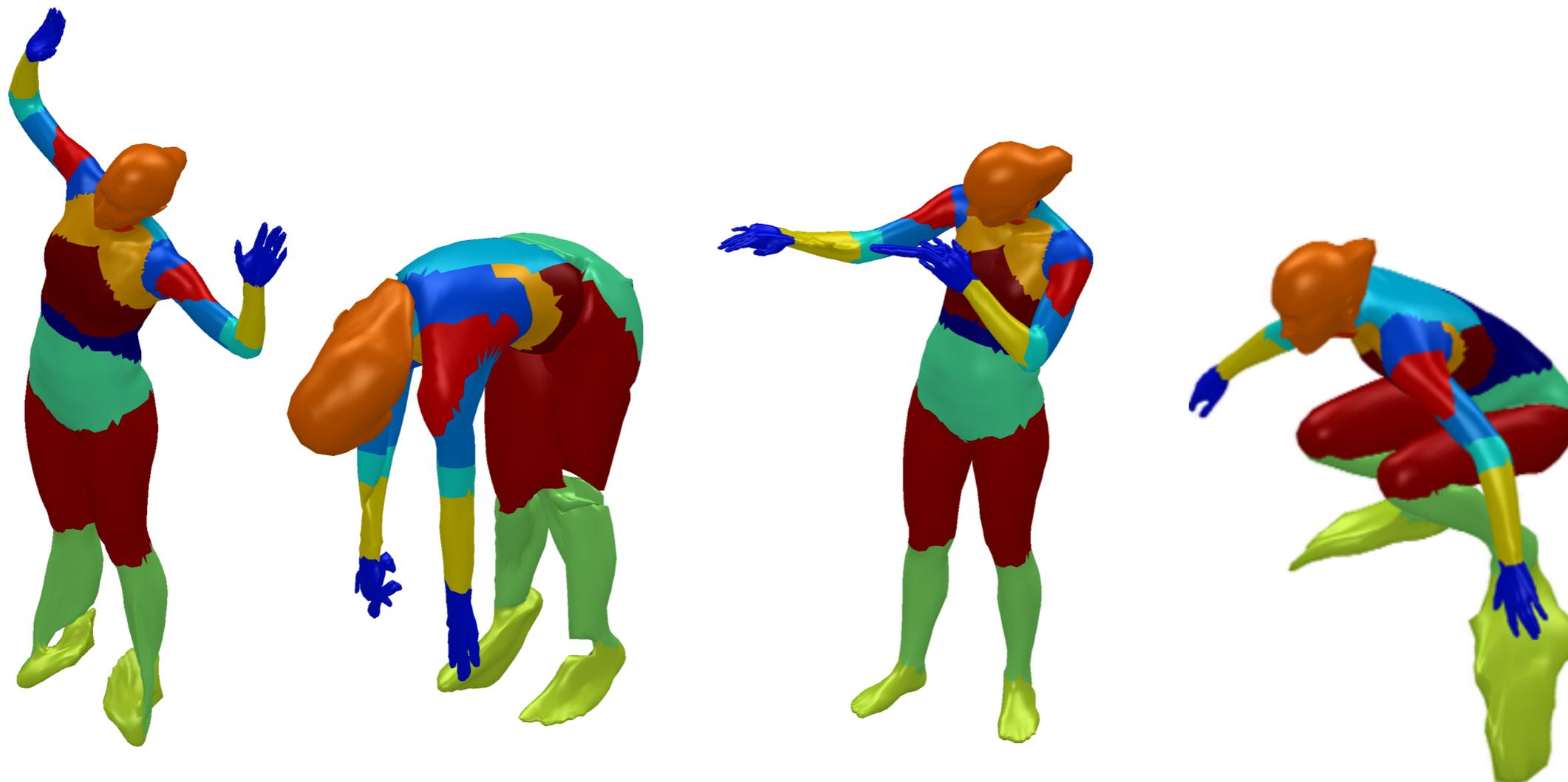
MoCap Likelihoods

$$\Sigma_{z_{gi}} \mid n_0, S_0 \sim \text{IW}(n_0, S_0),$$

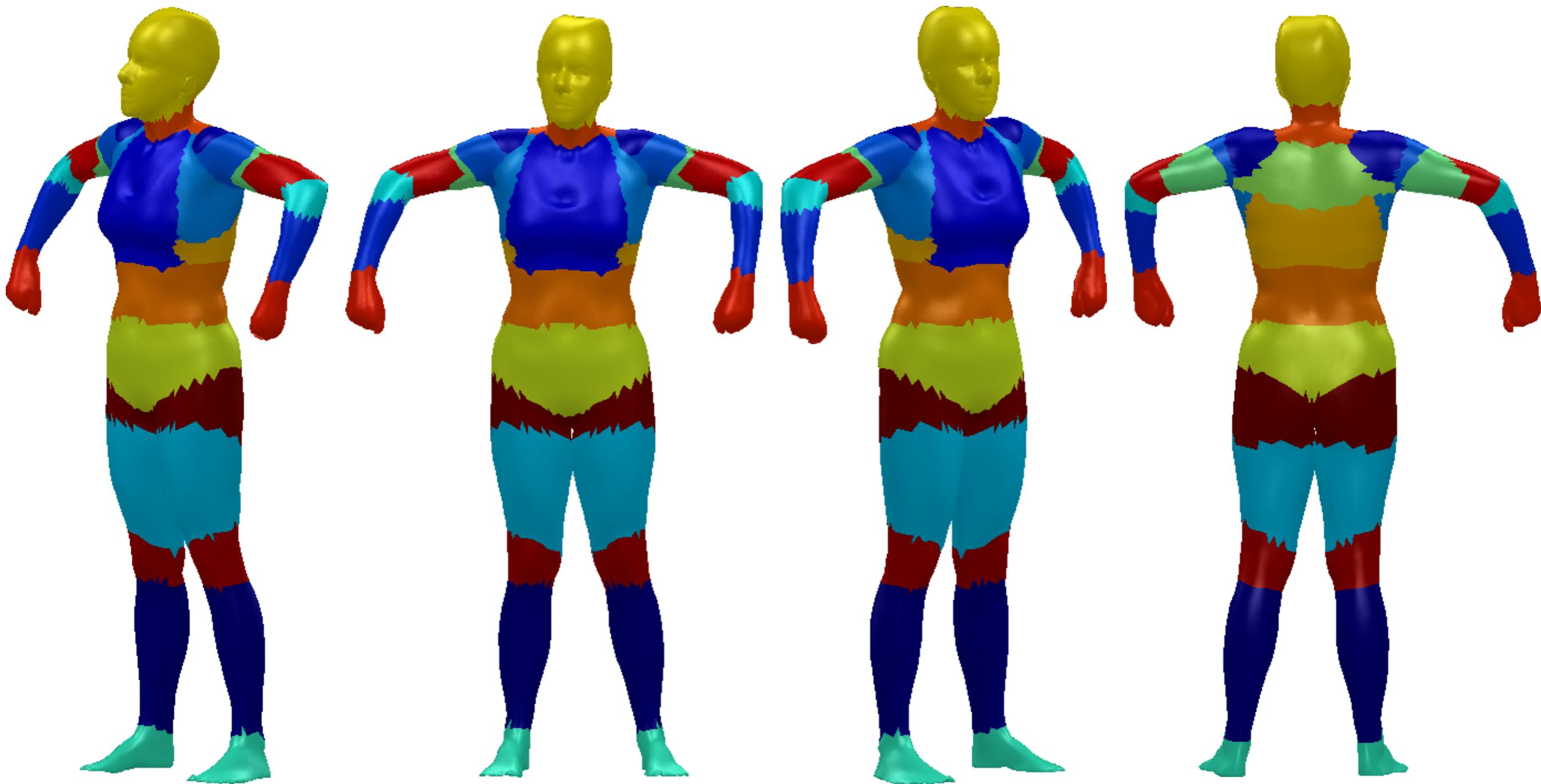
$$B_{z_{gi}} \mid M, \Sigma_{z_{gi}}, L \sim \mathcal{MN}(M, \Sigma_{z_{gi}}, L),$$

$$\epsilon_{z_{gi}} \sim \mathcal{N}(0, \Sigma_{z_{gi}}),$$

Moderate robustness to alignment errors



Inferred Segmentation



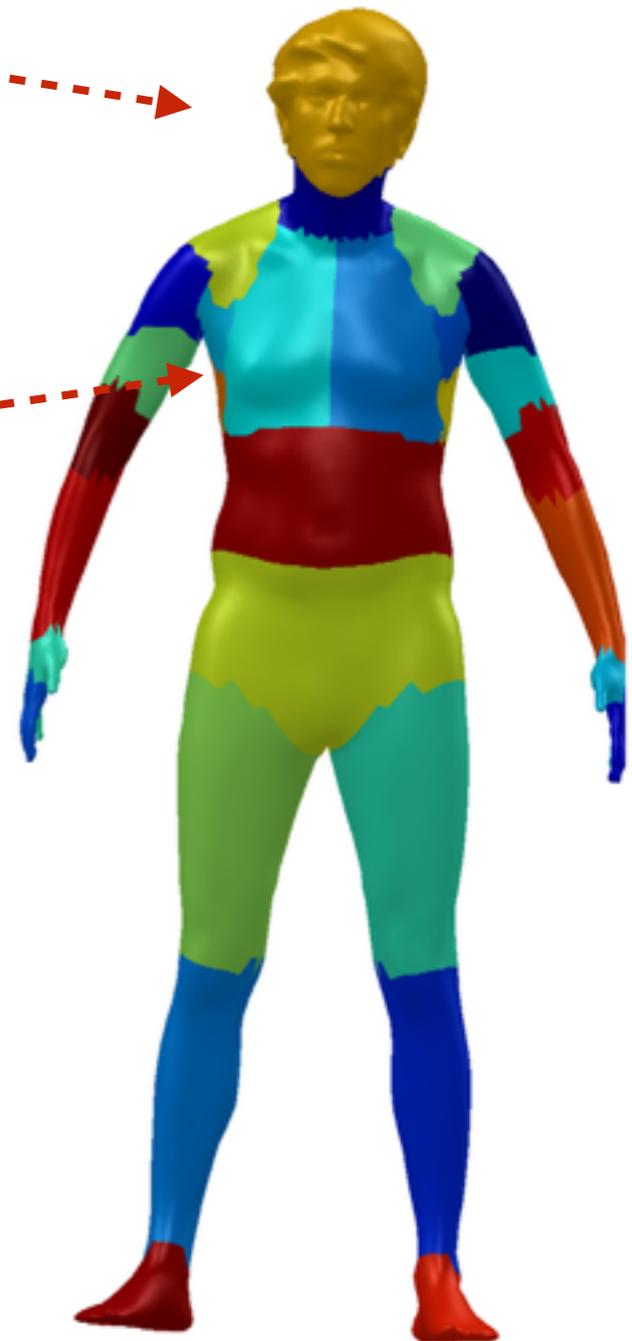
Segmentation with 20 Parts

Axial Symmetry

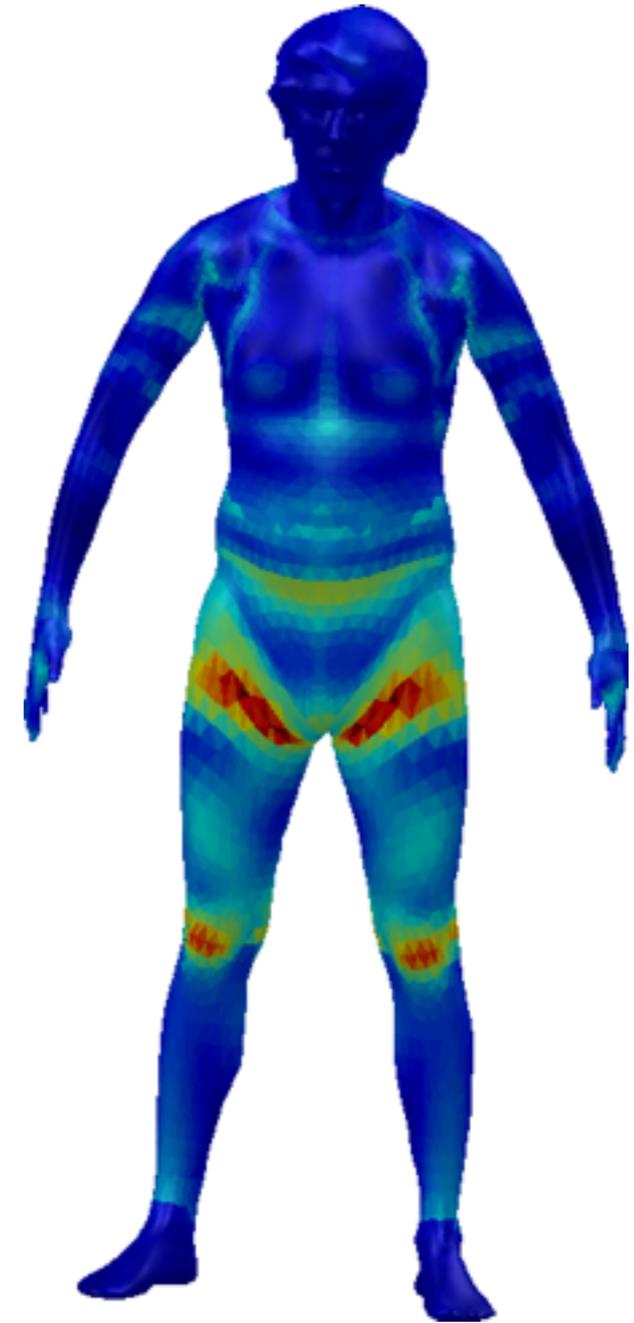
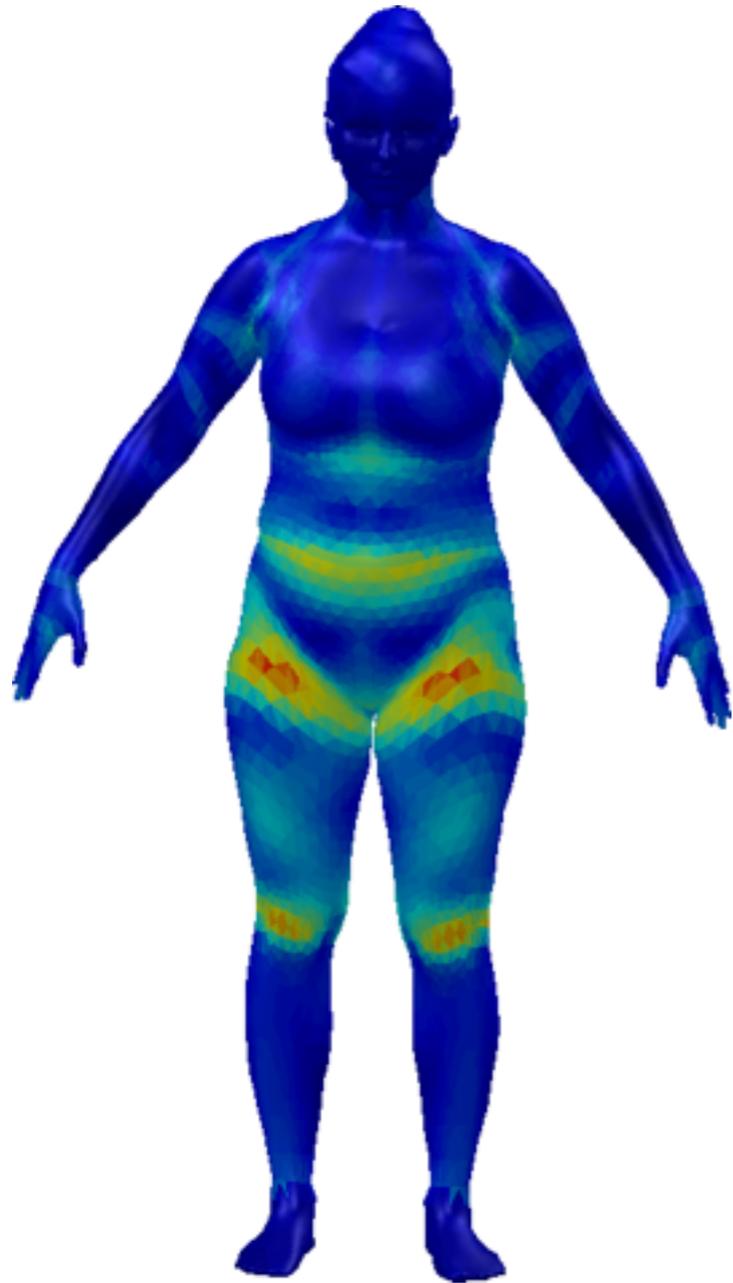


$$\frac{p(Y_{left}^{head} \cup Y_{right}^{head} | X_{left}^{head} \cup X_{right}^{head})}{p(Y_{left}^{head} | X_{left}^{head})p(Y_{right}^{head} | X_{right}^{head})} > 1$$

$$\frac{p(Y_{left}^{chest} \cup Y_{right}^{head} | X_{left}^{chest} \cup X_{right}^{chest})}{p(Y_{left}^{chest} | X_{left}^{chest})p(Y_{chest} | X_{right}^{chest})} < 1$$



*Only merge similarly moving parts
across axis of symmetry*



Measure of Rigidity

Inference

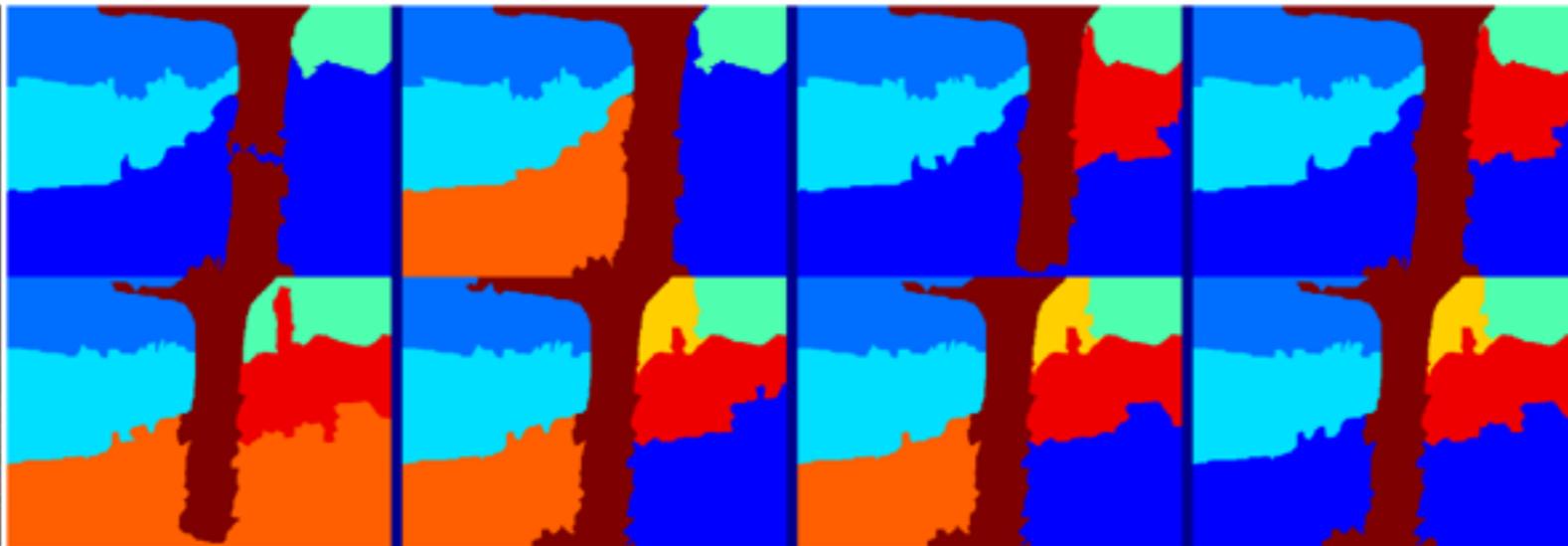
Algorithm 1: Hierarchical ddCRP sampler

For data instance $i \in \{1 \dots N_G\}$ jointly propose data and affected cluster links $\{\mathbf{c}^*, \mathbf{k}^*\} \leftarrow \text{ProposeLinks}(\mathbf{x}, \mathbf{k}, \mathbf{c}, \alpha_{1:G}, A^{1:G}, \alpha_0, A^0(\mathbf{c}))$.

Evaluate the proposal according to the Metropolis Hastings acceptance probability $a(\{\mathbf{c}^*, \mathbf{k}^*\}, \{\mathbf{c}, \mathbf{k}\})$. If the proposal is accepted, $\{\mathbf{c}^*, \mathbf{k}^*\}$ becomes the next state. If the proposal is rejected, the original configuration is retained.

For clusters $t \in T(\mathbf{c})$ resample cluster links via a Gibbs update:

$$k_t \sim p(k_t \mid \mathbf{k}_{-t}, \mathbf{c}, \mathbf{x}, \alpha_0, A^0(\mathbf{c})).$$



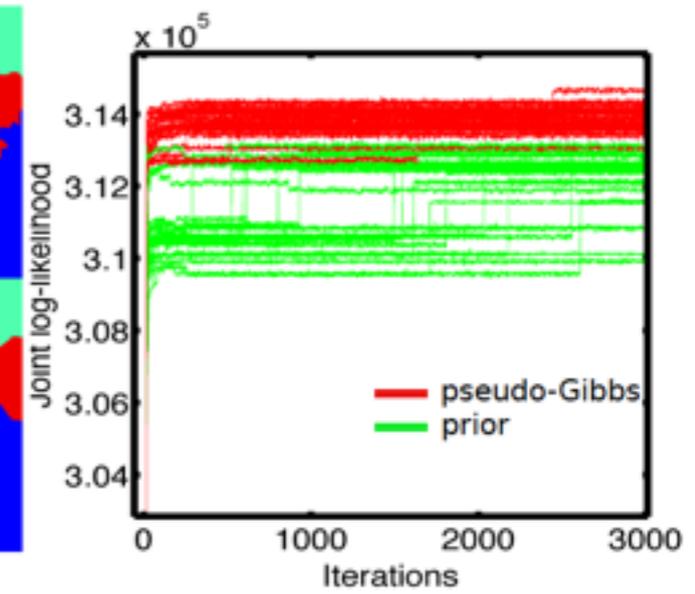
Worst prior proposal

Worst pseudo-Gibbs prop

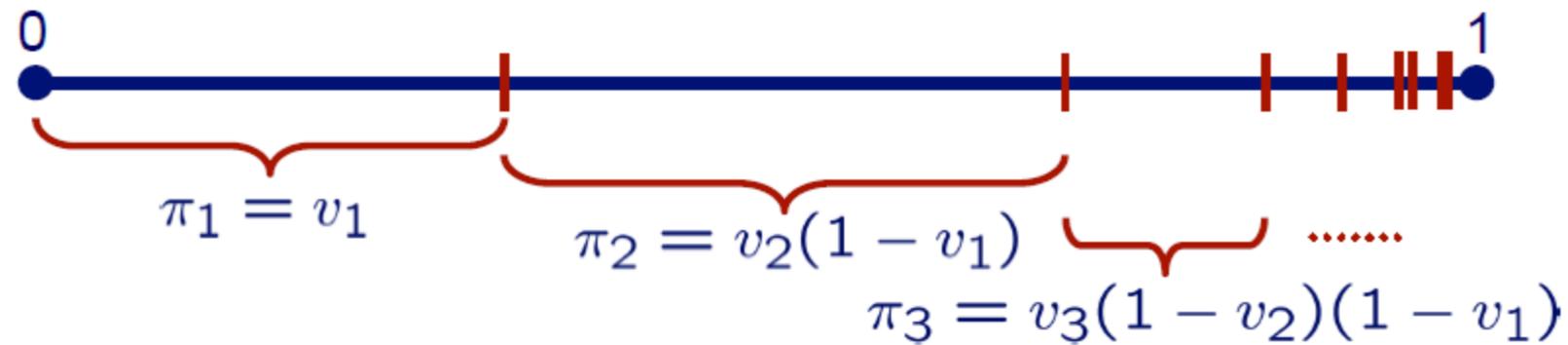
Best prior proposal

Best pseudo-Gibbs prop

Increasing Probability



Stick Breaking to Layers



$$v_k = \mathbb{P}(z_i = k \mid z_i \neq k - 1, \dots, 1)$$

Sequential Binary Sampler:

$$b_{ki} \sim \text{Bernoulli}(v_k)$$

$$z_i = \min\{k \mid b_{ki} = 1\}$$

- For each data instance i , go through the bins in order 1 through infinity.
- Toss a biased coin (with the probability of heads = v_k) for each bin .
- Pick the bin if the coin turns up heads

MCMC Learning

- Marginalize over the exponentially large space of latent links
 - MCMC samples
- Explore the marginal posterior of the auxiliary training model:

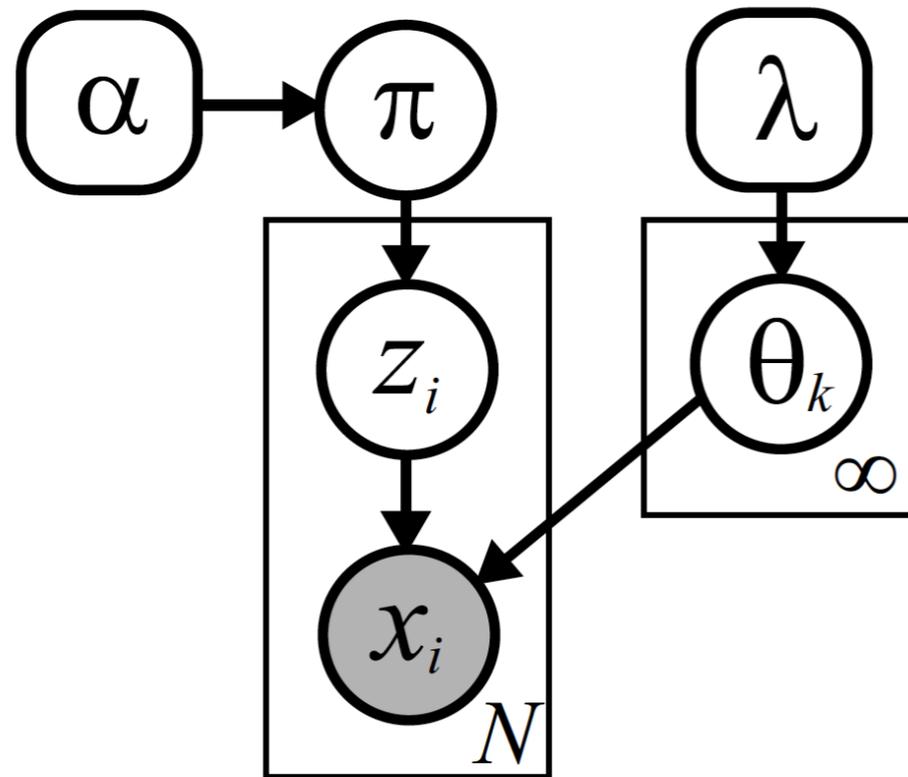
$$p(w_c | Y) = \sum_{\mathbf{c}} p(w_c, \mathbf{c} | Y) \approx \sum_{\mathbf{c}^{(s')}} p(w_c^{(s)}, \mathbf{c}^{(s')} | Y)$$

$$w_c^s, \mathbf{c}^s \sim p(w_c, \mathbf{c} | Y)$$

Random walk Proposal: $w_c^{t+1} \sim \mathcal{N}(w_c^{t+1} | w_c^t, \text{scale} \times \mathbf{I})$

Gibbs Step: $c_{di} | \mathbf{c}_{-di}, w_c^*, Y \sim p(c_{di} | w_c^*) \delta(z(\mathbf{c}_d), y_d)$

Bayesian Nonparametric Priors



$$G(\theta) = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}(\theta)$$

$$\sum_{k=1}^{\infty} \pi_k = 1 \quad 0 \leq \pi_k \leq 1$$

$$\theta_k \sim H(\lambda)$$

Pitman-Yor Process

Power Law Behavior

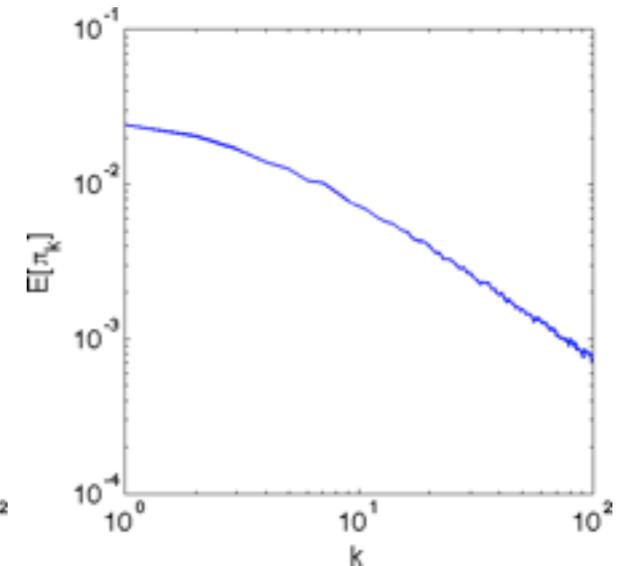
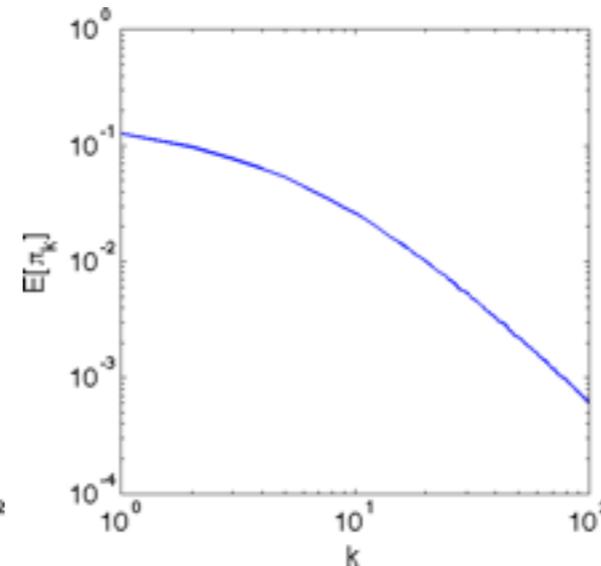
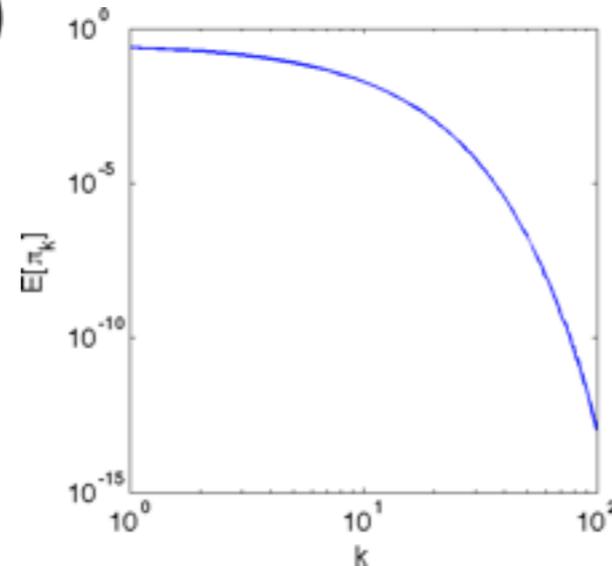
$$E[w_k] = \frac{1 - \alpha_a}{(1 + \alpha_b + (k - 1)\alpha_a)}$$

$$\pi_k = w_k \prod_{l=1}^{k-1} (1 - w_l)$$

$$w_k \sim \text{Beta}(a_k, b_k)$$

$$a_k = 1 - \alpha_a$$

$$b_k = \alpha_b + k\alpha_a$$



Number of unique clusters in N observations: $O(\alpha_b N^{\alpha_a})$

Expected size of sorted component k : $O(k^{-\frac{1}{\alpha_a}})$

Hierarchical ddCRP

Sample local links:

$$p(c_{gi} = gj | \alpha_g, A^g) \propto \begin{cases} A_{ij}^g & i \neq j, \\ \alpha_g & i = j. \end{cases}$$

$$\Lambda_g = z(\mathbf{c}_g)$$

Sample global links:

$$p(k_t = s | \alpha_0, A^0(\mathbf{c})) \propto \begin{cases} A_{ts}^0(\mathbf{c}) & t \neq s, \\ \alpha_0 & t = s. \end{cases}$$

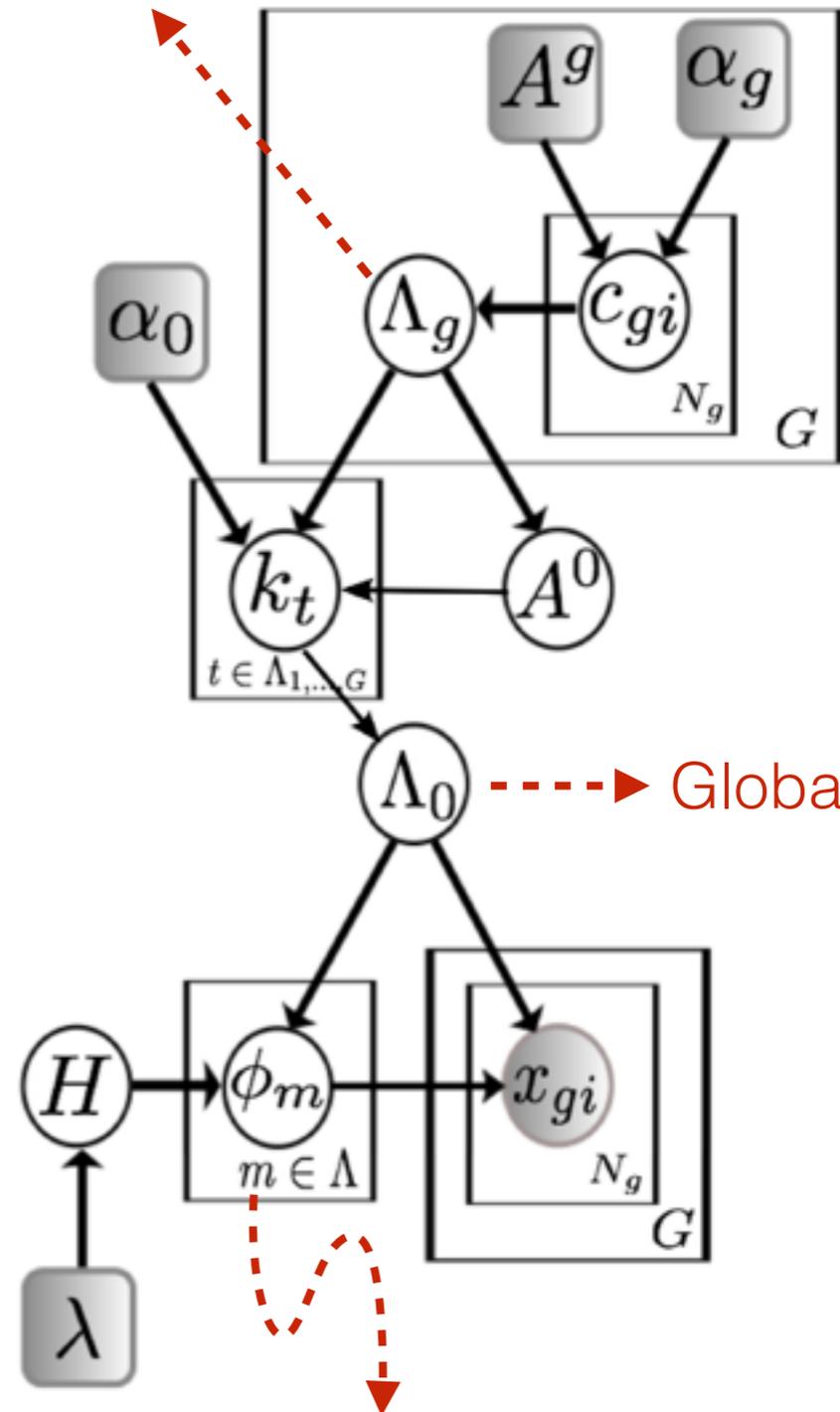
$$\Lambda_0 = z(\mathbf{k})$$

Sample data generating parameters:

$$\phi_m \sim H(\lambda), \forall m \in \Lambda_0$$

$$x_i \sim \phi_m, \forall i \in m$$

Group Specific Partitions



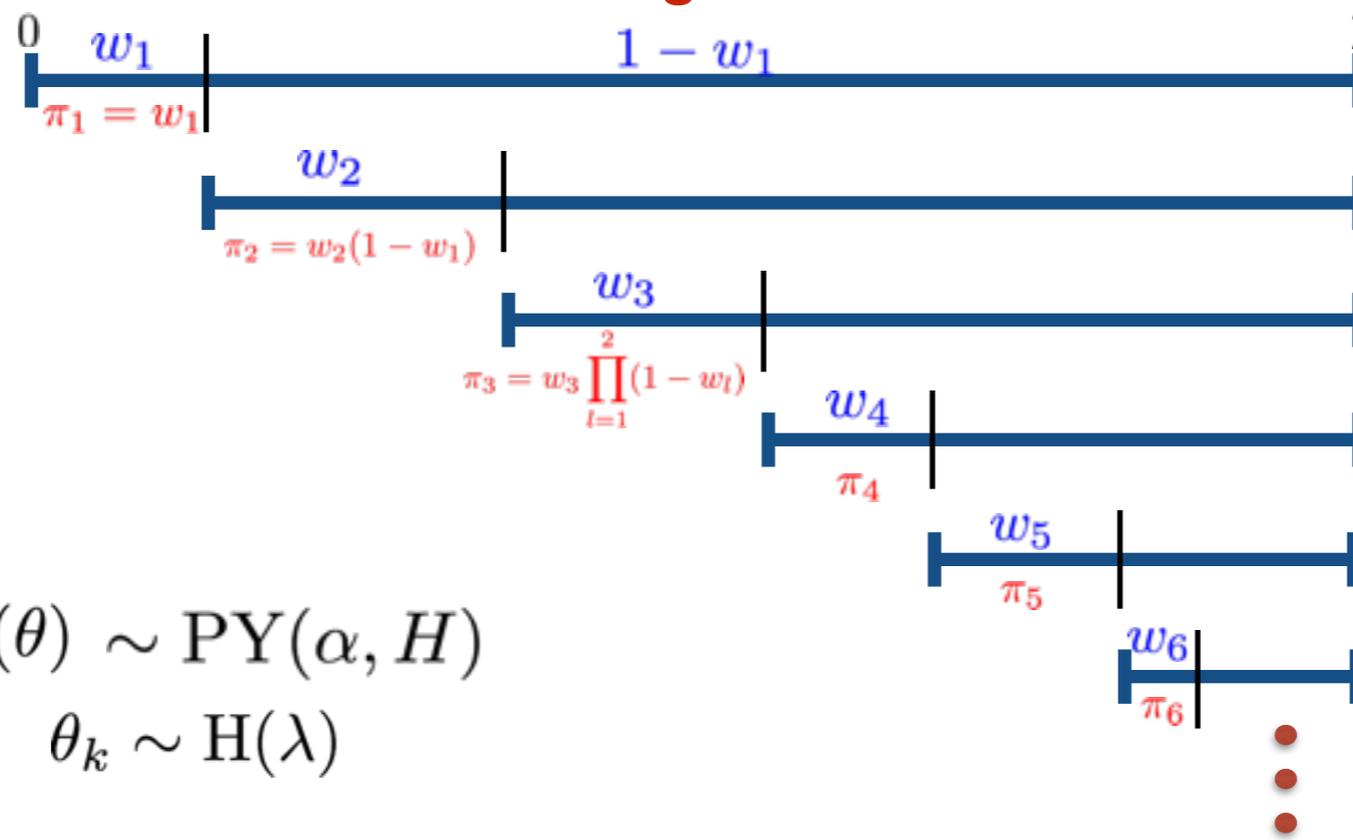
Components shared across groups

Pitman-Yor Process

- The Pitman-Yor process defines a distribution on infinite discrete measures, or partitions

$$\pi_k = w_k \prod_{l=1}^{k-1} (1 - w_l) \quad w_k \sim \text{Beta}(1 - \alpha_a, \alpha_b + k\alpha_a)$$

Stick Breaking Construction:



$$G(\theta) = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}(\theta) \sim \text{PY}(\alpha, H)$$

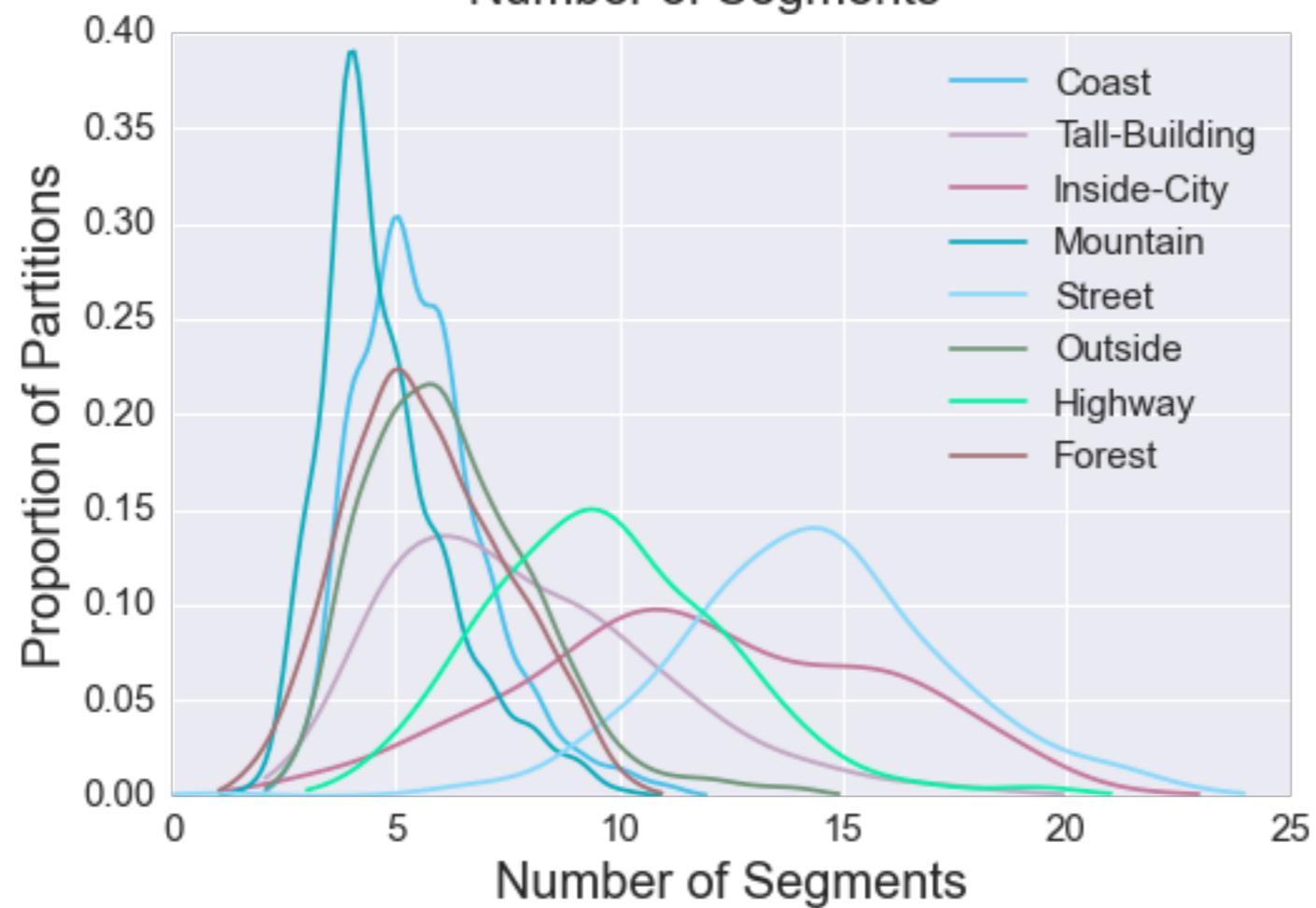
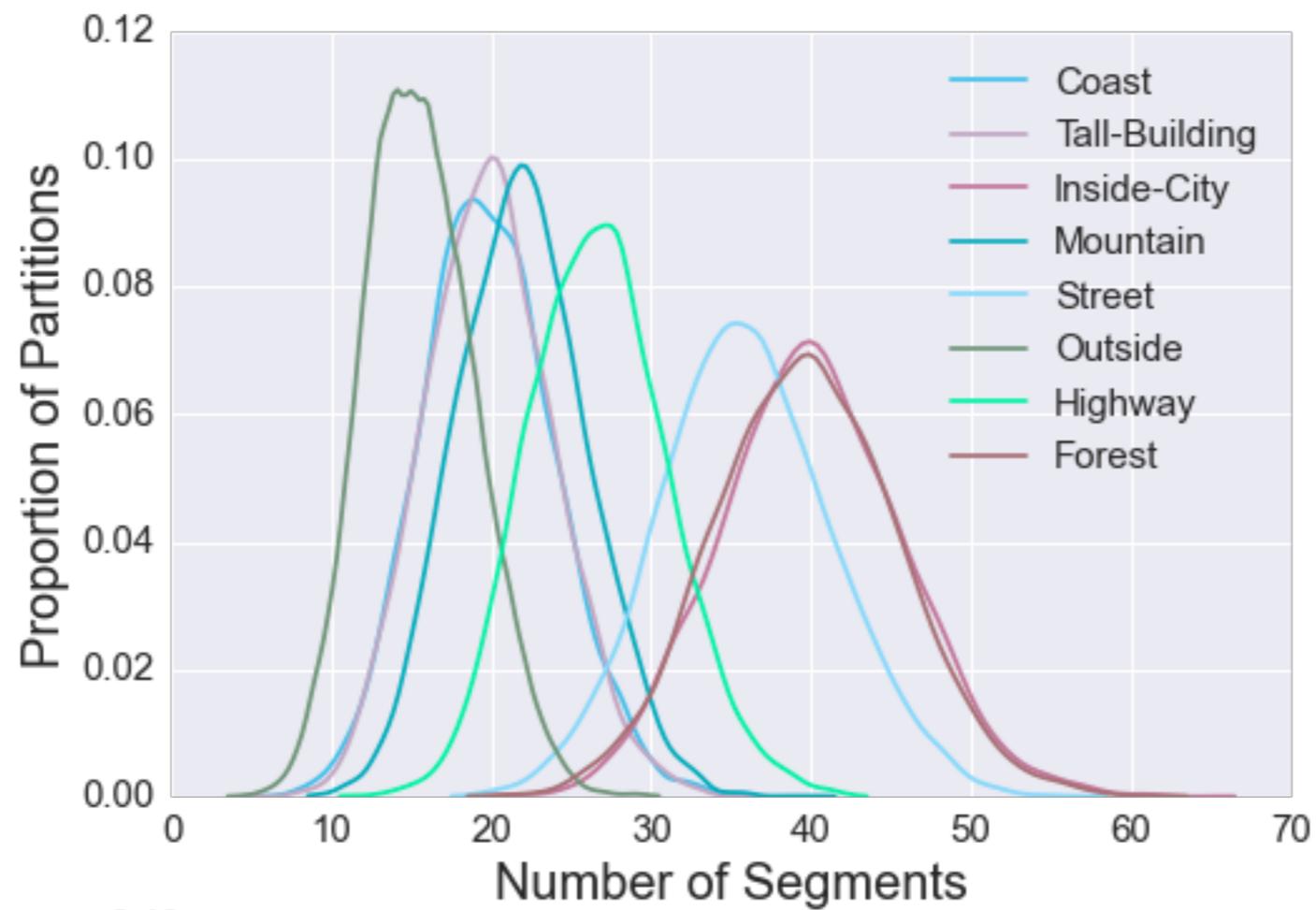
$$\theta_k \sim H(\lambda)$$

Sethuraman, 1994
Ishwaran and James,
2001

Video Segmentation

$$P = \frac{\sum_{i=1}^M \left[\left\{ \sum_{s \in \mathcal{S}} \max_{g \in \mathbb{G}_i} |s \cap g| \right\} - \max_{g \in \mathbb{G}_i} |g| \right]}{M|\mathcal{S}| - \sum_{i=1}^M \max_{g \in \mathbb{G}_i} |g|}$$
$$R = \frac{\sum_{i=1}^M \sum_{g \in \mathbb{G}_i} \{ \max_{s \in \mathcal{S}} |s \cap g| - 1 \}}{\sum_{i=1}^M \{ |\mathbb{G}_i| - \Gamma_{\mathbb{G}_i} \}}$$

VPR



Approximate Bayesian Computation

Algorithm 3 Likelihood-free MCMC sampler

Use Algorithm 2 to get a realisation $(\boldsymbol{\theta}^{(0)}, \mathbf{z}^{(0)})$ from the ABC target distribution $\pi_\varepsilon(\boldsymbol{\theta}, \mathbf{z}|\mathbf{y})$

for $t = 1$ to N **do**

 Generate $\boldsymbol{\theta}'$ from the Markov kernel $q(\cdot|\boldsymbol{\theta}^{(t-1)})$,

 Generate \mathbf{z}' from the likelihood $f(\cdot|\boldsymbol{\theta}')$,

 Generate u from $\mathcal{U}_{[0,1]}$,

if $u \leq \frac{\pi(\boldsymbol{\theta}')q(\boldsymbol{\theta}^{(t-1)}|\boldsymbol{\theta}')}{\pi(\boldsymbol{\theta}^{(t-1)})q(\boldsymbol{\theta}'|\boldsymbol{\theta}^{(t-1)})}$ and $\rho\{\eta(\mathbf{z}'), \eta(\mathbf{y})\} \leq \varepsilon$ **then**

 set $(\boldsymbol{\theta}^{(t)}, \mathbf{z}^{(t)}) = (\boldsymbol{\theta}', \mathbf{z}')$

else

$(\boldsymbol{\theta}^{(t)}, \mathbf{z}^{(t)}) = (\boldsymbol{\theta}^{(t-1)}, \mathbf{z}^{(t-1)})$,

end if

end for
