

# Fast search for Dirichlet process mixture models

**Hal Daumé III**

School of Computing  
University of Utah

me@hal3.name



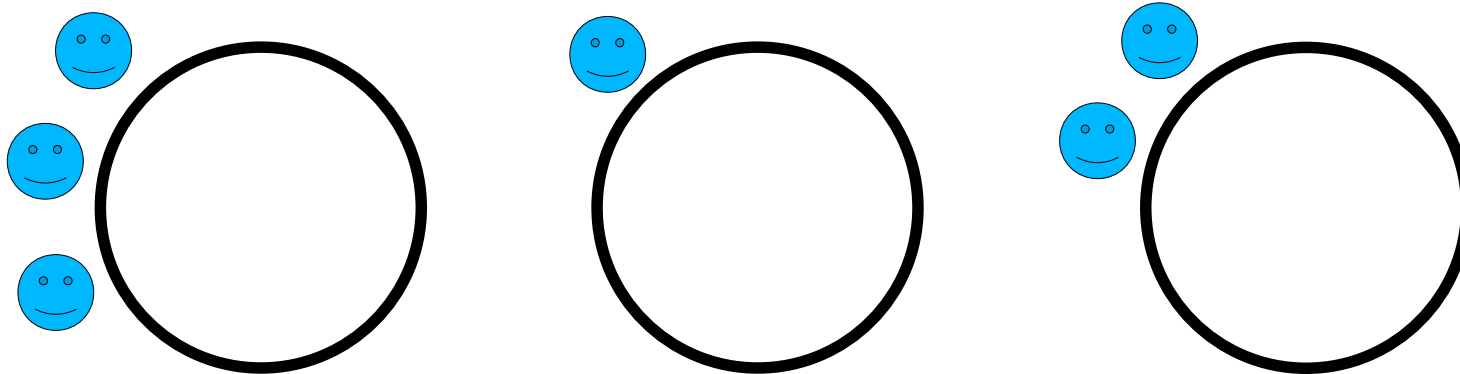
# Dirichlet Process Mixture Models

- Non-parametric Bayesian density estimation
- Frequently used to solve clustering problem: choose “K”
- Applications:
  - vision
  - data mining
  - computational biology

## Personal Observation:

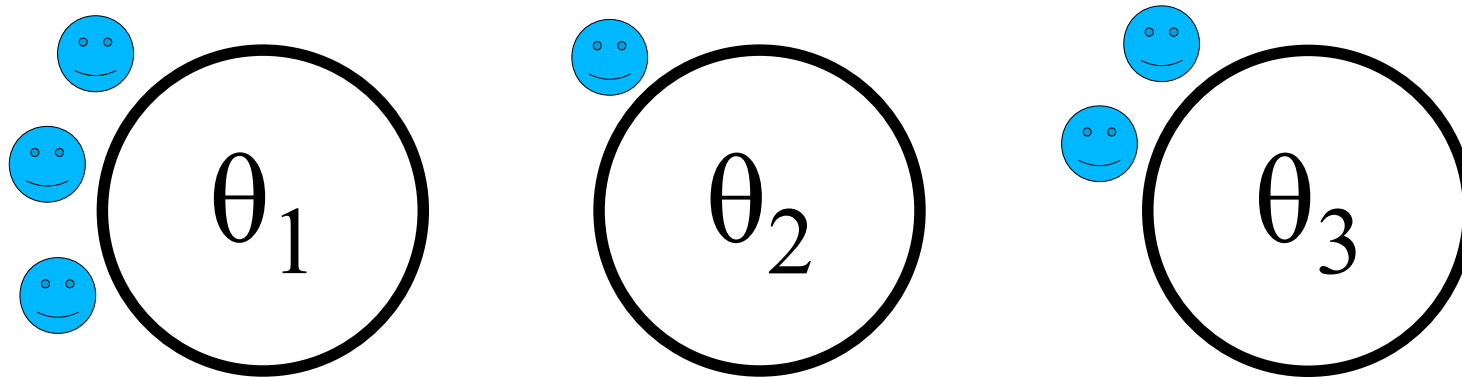
Samplers slow on huge data sets (10k+ elements)  
Very sensitive to initialization

# Chinese Restaurant Process



- Customers enter a restaurant sequentially
- The  $M$ th customer chooses a table by:
  - Sit at table with  $N$  customers with probability  $N/(\alpha+M-1)$
  - Sit at unoccupied table with probability  $\alpha/(\alpha+M-1)$

# Dirichlet Process Mixture Models



- Data point = customer
- Cluster = table
- Each table gets a parameter
- Data points are generated according to a likelihood  $F$

$$p(X | c) = \int d\theta_{1:K} \left[ \prod_k G_0(\theta_k) \right] \left[ \prod_n F(x_n | \theta_{c_n}) \right]$$

$c_n$  = table of  $n$ th customer

# Inference Summary

- Run MCMC sampler for a bunch of iterations
  - Use different initialization
- From set of samples, choose one with highest posterior probability

If all we want is the highest probability assignment, why not just try to find it directly?

*(If you really want to be Bayesian, use this assignment to initialize sampling)*

# Ordered Search

**Input:** data, beam size, scoring function

**Output:** clustering

initialize  $Q$ , a max-queue of partial clusterings

while  $Q$  is not empty

    remove a partial cluster  $c$  from  $Q$

    if  $c$  covers all the data, return it

    try extending  $c$  by a single data point

    put all  $K+1$  options into  $Q$  with scores

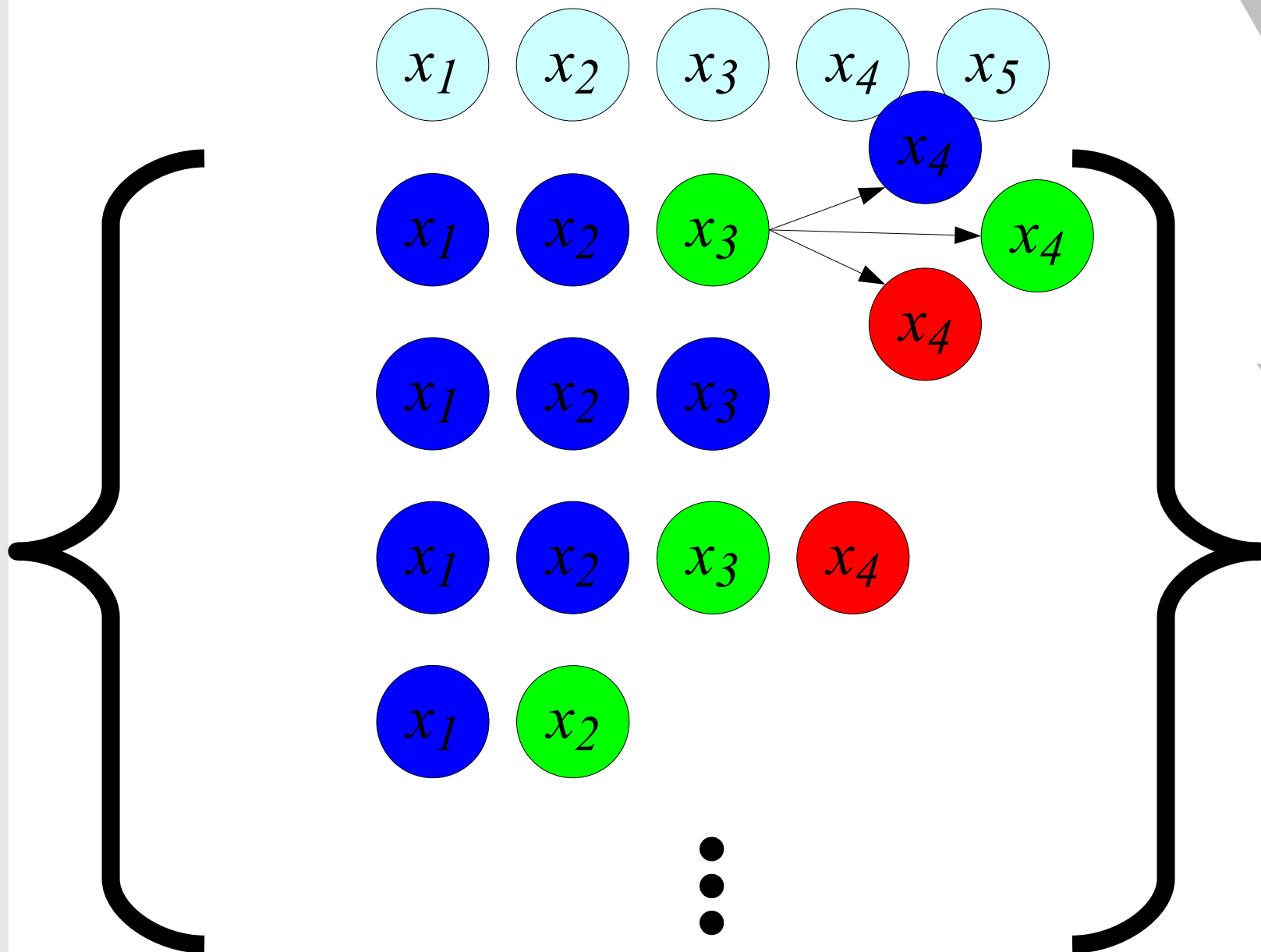
    if  $|Q| >$  beam size, drop elements

**Optimal, if:**

    beam size = infinity

    scoring function *overestimates* true best probability

# Ordered Search in Pictures



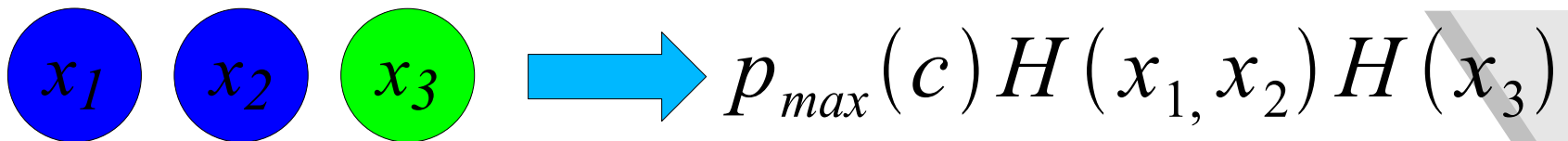
# Trivial Scoring Function

- Only account for already-clustered data:

$$g^{Triv}(c, x) = p_{max}(c) \prod_{k \in c} H(\{x_{c=k}\})$$

$$H(X) = \int d\theta G_0(\theta) \prod_{x \in X} F(x | \theta)$$

- $p_{max}(c)$  can be computed exactly



$$x_1 \quad x_2 \quad x_3 \quad \longrightarrow \quad p_{max}(c) H(x_1, x_2) H(x_3)$$



# Tighter Scoring Function

- Use trivial score for already-clustered data
- Approximate optimal score for future data:
  - For each data point, put in existing or new cluster
  - Then, conditioned on that choice, cluster remaining
    - Assume each remaining point is optimally placed

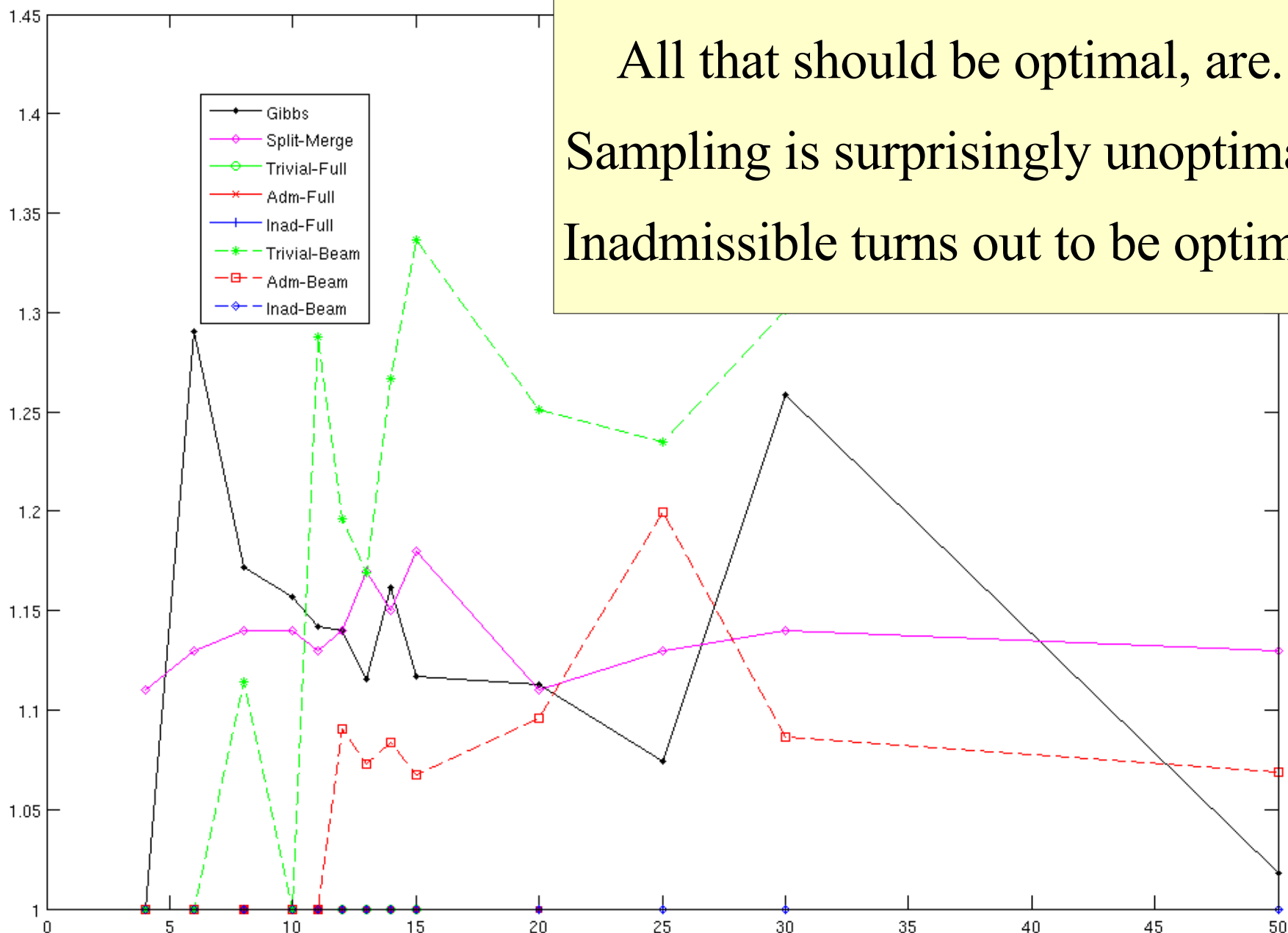
# An Inadmissible Scoring Function

- Just use marginals for unclustered points:

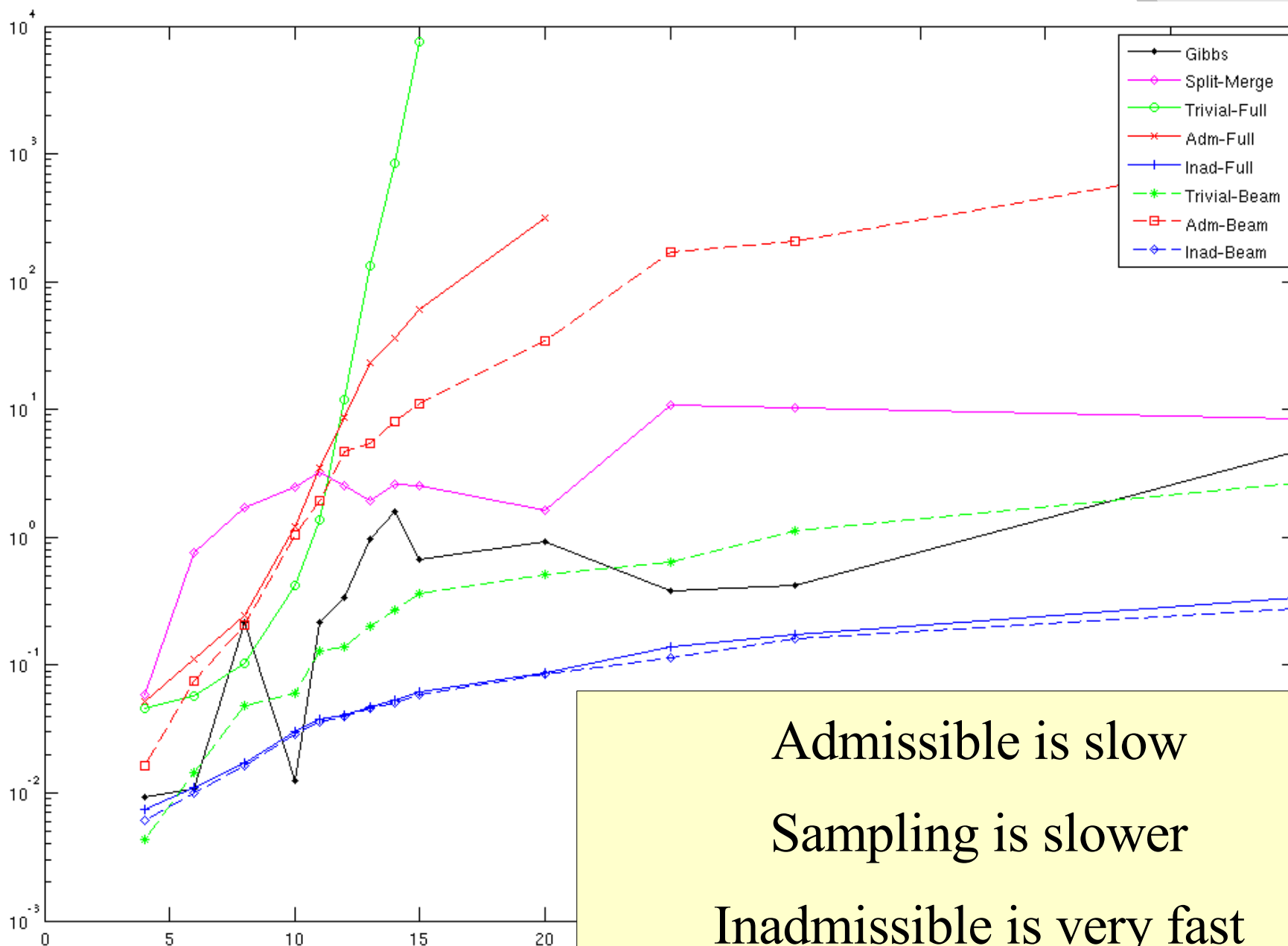
$$g^{Inad}(c, x) = g^{Triv}(c, x) \prod_{n=|c|+1}^N H(x_n)$$

- Inadmissible because H is not monotonic in conditioning (even for exponential family)

# Artificial Data: Likelihoods Ratio

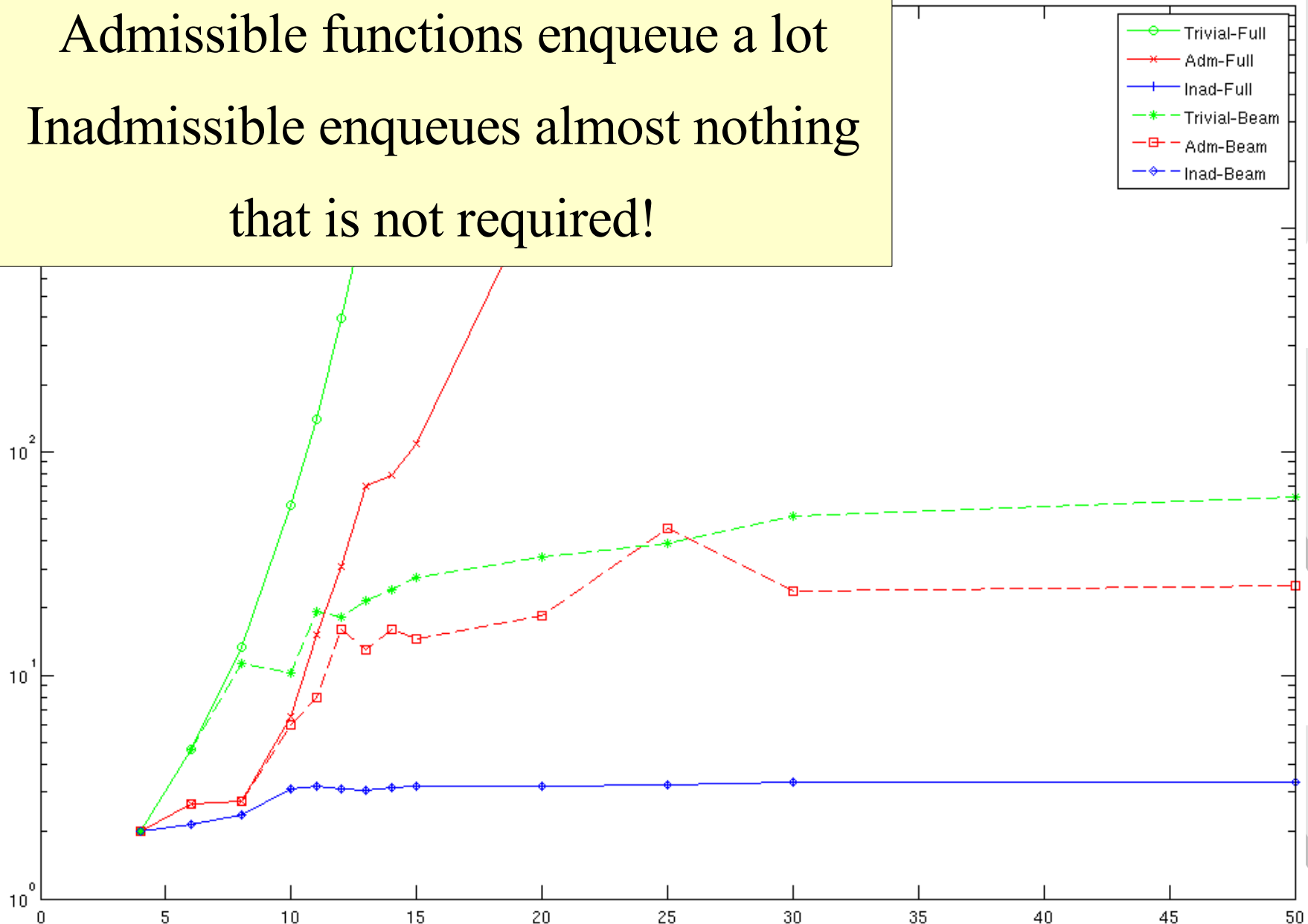


# Artificial Data: Speed (Seconds)



# Artificial Data: # of Enqueued Points

Admissible functions enqueue a lot  
Inadmissible enqueues almost nothing  
that is not required!



# Real Data: MNIST

- Handwritten numbers 0-9, 28x28 pixels
- Preprocess with PCA to 50 dimensions
- Run on: 3000, 12,000 and 60,000 images
- Use inadmissible heuristic with (large) 100 beam

	<u>3k</u>	<u>12k</u>	<u>60k</u>
<b>Search</b>	11s <i>2.04e5</i>	105s <i>8.02e5</i>	15m <i>3.96e6</i>
<b>Gibbs</b>	40s/i <i>2.09e5</i>	18m/i <i>8.34e5</i>	7h/i <i>4.2e6</i>
<b>S-M</b>	85s/i <i>2.05e5</i>	35m/i <i>8.15e5</i>	12h/i <i>4.1e6</i>

# Real Data: NIPS Papers

- NIPS 1-12
- 1740 documents, vocabulary of 13k words
  - Drop top 10, retain remaining top 1k
- Conjugate Dirichlet/Multinomial DP
- Order examples by increasing marginal likelihood

<b>Search</b>	$2.441e6$ (32s) $2.474e6$ ( <i>reverse order</i> ) $2.449e6$ ( <i>random order</i> )
---------------	--

<b>Gibbs</b>	$3.2e6$ (1h)
--------------	--------------

<b>S-M</b>	$3.0e6$ (1.5h)
------------	----------------

# Discussion

**Sampling often fails to find MAP**

**Search can do much better**

**Limited to conjugate distributions**

**Cannot re-estimate hyperparameters**

**Can cluster 270 images / second in matlab**

**Further acceleration possible with  
clever data structures**

**Thanks! Questions?**

code at <http://hal3.name/DPsearch>



# Inference I – Gibbs Sampling

## Collapsed Gibbs sampler:

- Initialize clusters
- For a number of iterations:
  - Assign each data point  $x_n$  to cluster  $c_k$  with probability:

$$\frac{N_k}{\alpha + N - 1} \int d\theta G_0(\theta) F(x_n|\theta) \prod_{m \in c_k} F(x_m|\theta)$$

or to a new cluster with probability

$$\frac{\alpha}{\alpha + N - 1} \int d\theta G_0(\theta) F(x_n|\theta)$$

$$= H(x_n | \{x_{c=k}\})$$

H is the posterior probability of  $x$ , conditioned on the set of  $x$  that fall into the proposed cluster

# Inference II – Metropolis-Hastings

## Collapsed Split-Merge sampler:

- Initialize clusters
- For a number of iterations:
  - Choose two data points  $x_n$  and  $x_m$  at random
    - If  $c_n = c_m$ , split this cluster with a Gibbs pass
    - otherwise, merge the two clusters
  - Then perform a collapsed Gibbs pass

# Versus Variational

