# Accelerated Gibbs Sampling for the Indian Buffet Process

Finale Doshi-Velez, Cambridge/MIT
Zoubin Ghahramani, Cambridge

# Motivation

Bilinear models of the form

$$X = UV + E$$
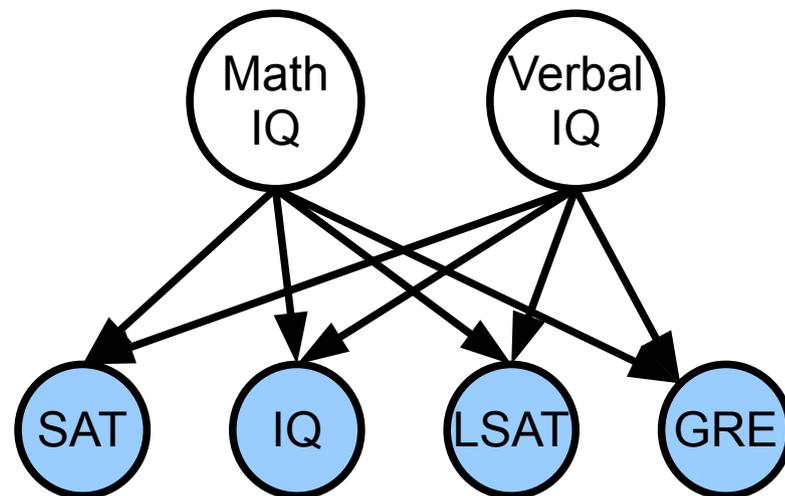
data = matrix product + error

are very common in machine learning.

UNIVERSITY OF
CAMBRIDGE

# Examples

## Factor Analysis
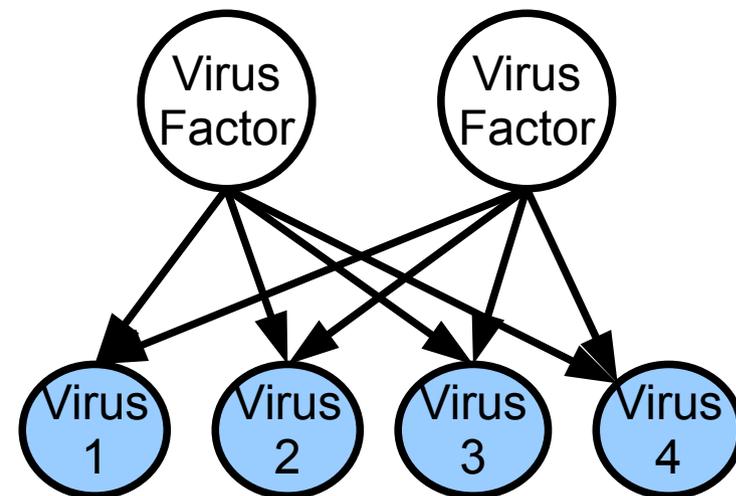
$$Y = LX + E$$

UNIVERSITY OF
CAMBRIDGE

# Examples

Factor Analysis

$$Y = LX + E$$

Probabilistic PCA

$$T = WX + E$$

# Examples

Factor Analysis

$$Y = LX + E$$

Probabilistic PCA

$$T = WX + E$$

Probabilistic Matrix Factorization

$$X = UV + E$$



User-Movie Ratings = User Features · Movie Features

UNIVERSITY OF CAMBRIDGE
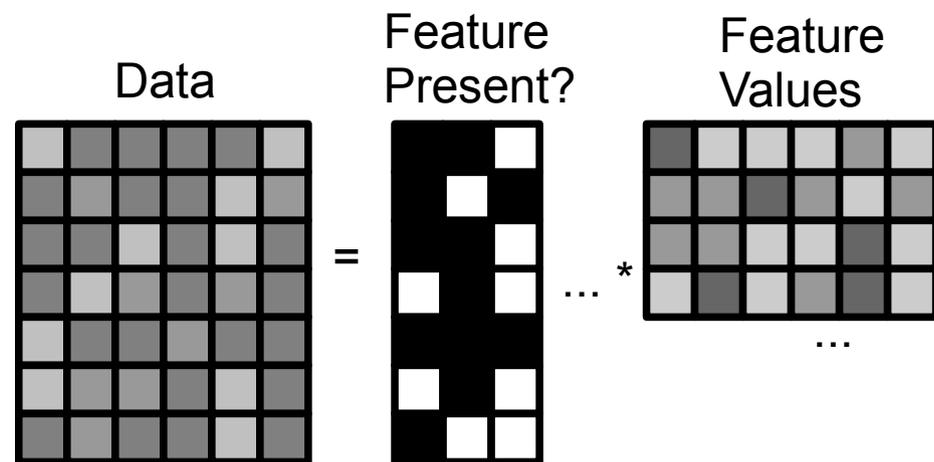
# Examples

Factor Analysis

$$Y = LX + E$$

Probabilistic PCA

$$T = WX + E$$

Probabilistic Matrix Factorization

$$X = UV + E$$

Indian Buffet Process with a linear likelihood

$$X = ZA + E$$



Data = Feature Present? ... * Feature Values ...

UNIVERSITY OF
CAMBRIDGE

# Motivation

- We are interested in doing large-scale Bayesian inference in these models (focus on the IBP for now):
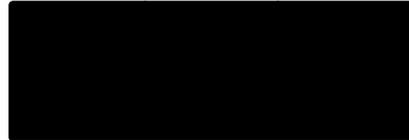
$$X = ZA + E$$

- Suppose

  - We <u>can</u> compute $P(X|Z)$ , <u>but</u> it's expensive
  - We <u>can</u> compute $P(A|X,Z)$
  - We <u>cannot</u> compute $P(Z,A|X)$

- We develop a fast sampler for inference in these models.

UNIVERSITY OF CAMBRIDGE

# Indian Buffet Process

Customers enter an "infinite buffet" one at a time and

- Sample a previously sampled dish based on its popularity.

- Sample Poisson( alpha / n ) new dishes.



...

UNIVERSITY OF CAMBRIDGE

# Indian Buffet Process

Customers enter an "infinite buffet" one at a time and

- Sample a previously sampled dish based on its popularity.

- Sample Poisson( alpha / n ) new dishes.

UNIVERSITY OF
CAMBRIDGE

# Indian Buffet Process
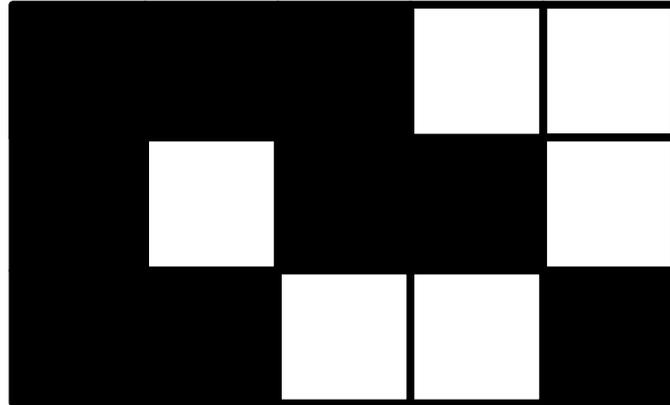
Customers enter an "infinite buffet" one at a time and

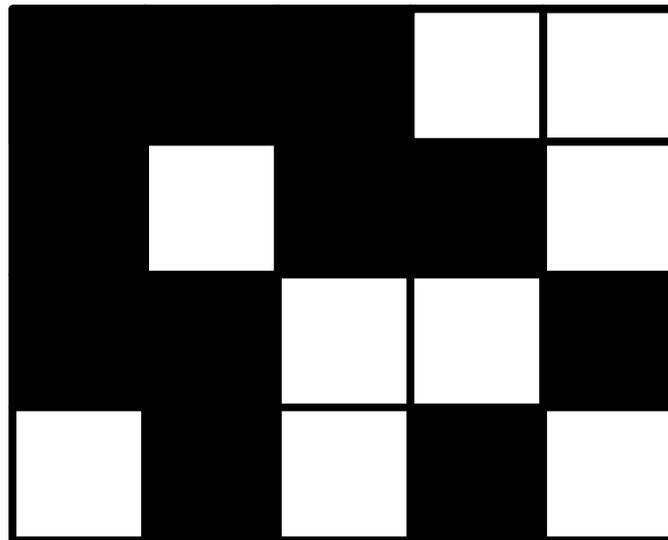- Sample a previously sampled dish based on its popularity.

- Sample Poisson( alpha / n ) new dishes.

# Indian Buffet Process
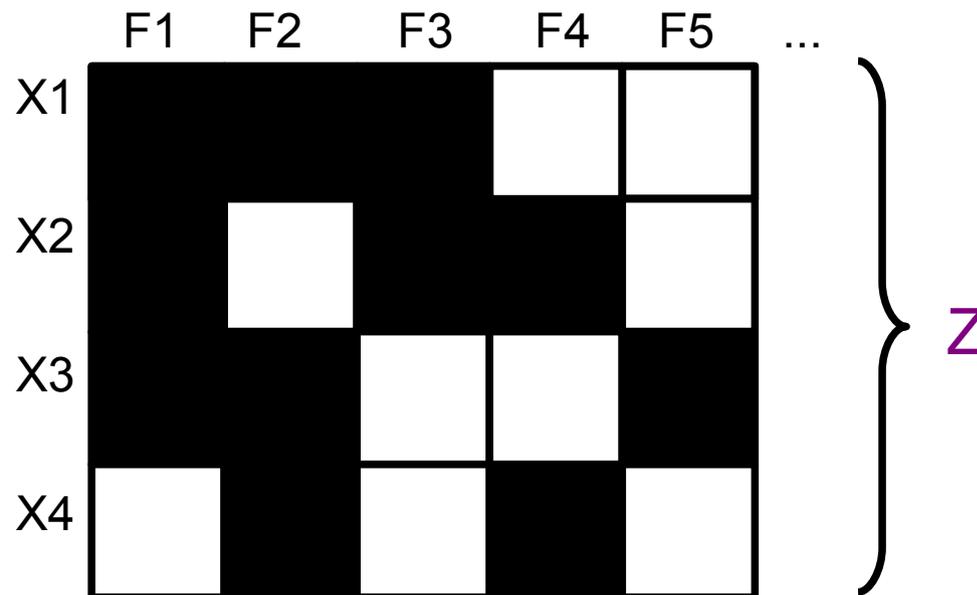
Customers enter an "infinite buffet" one at a time and

- Sample a previously sampled dish based on its popularity.

- Sample Poisson( alpha / n ) new dishes.

UNIVERSITY OF
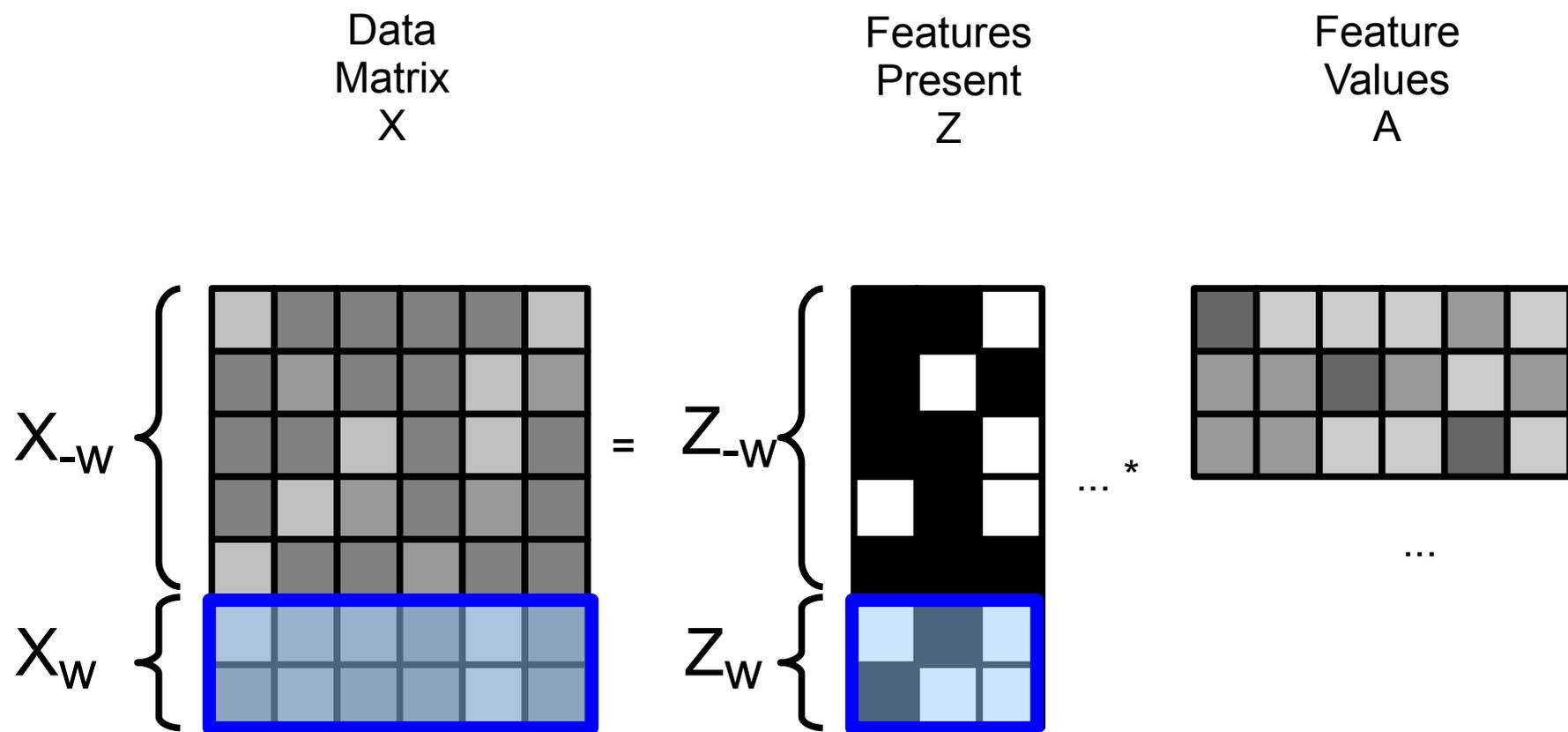CAMBRIDGE

# Indian Buffet Process

Result is a non-parametric prior on feature assignments—a general tool for many latent feature models—with some nice properties:

- Observations are exchangeable.

- Infinite features, but finite datasets contain a finite number of features.
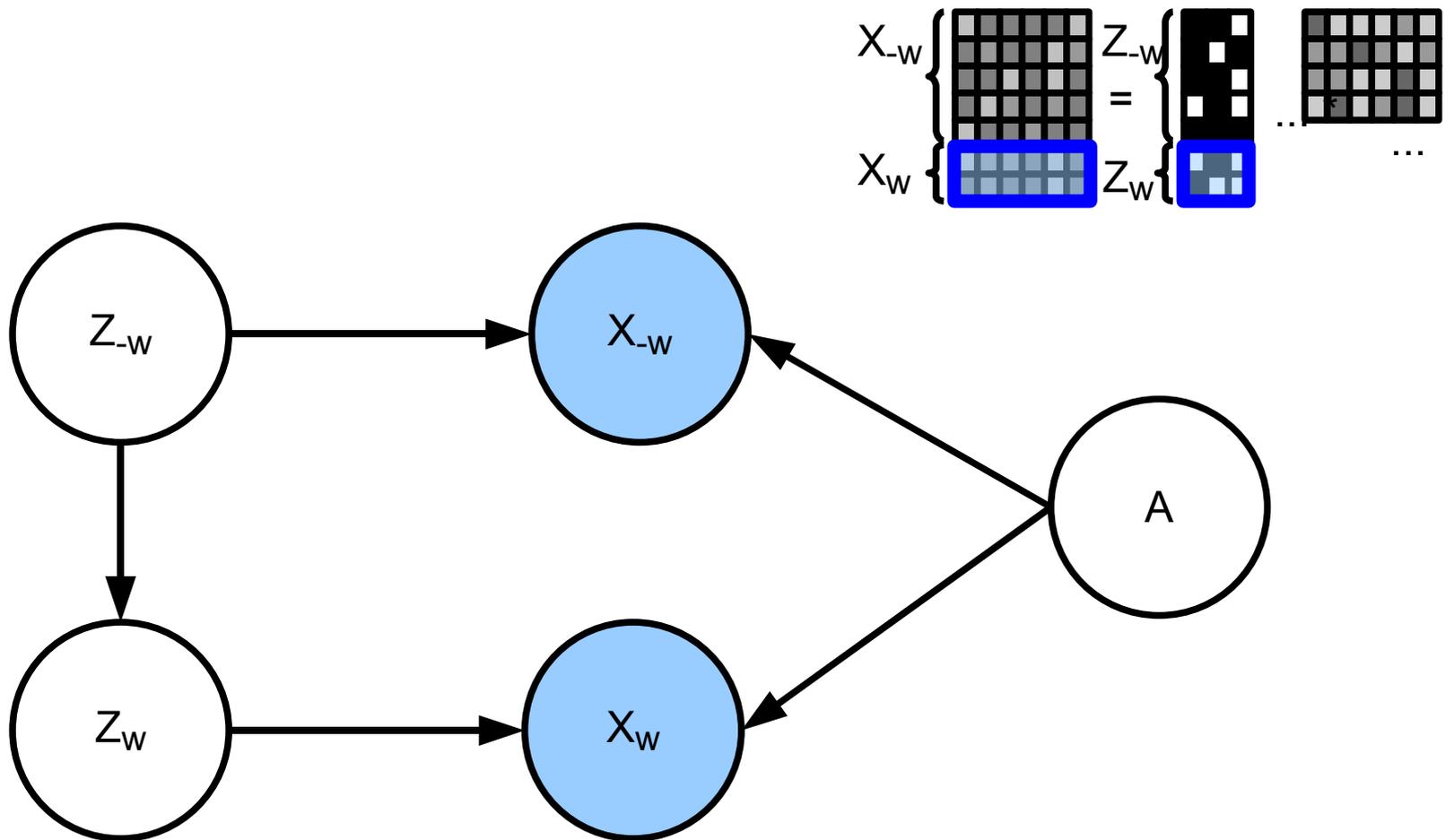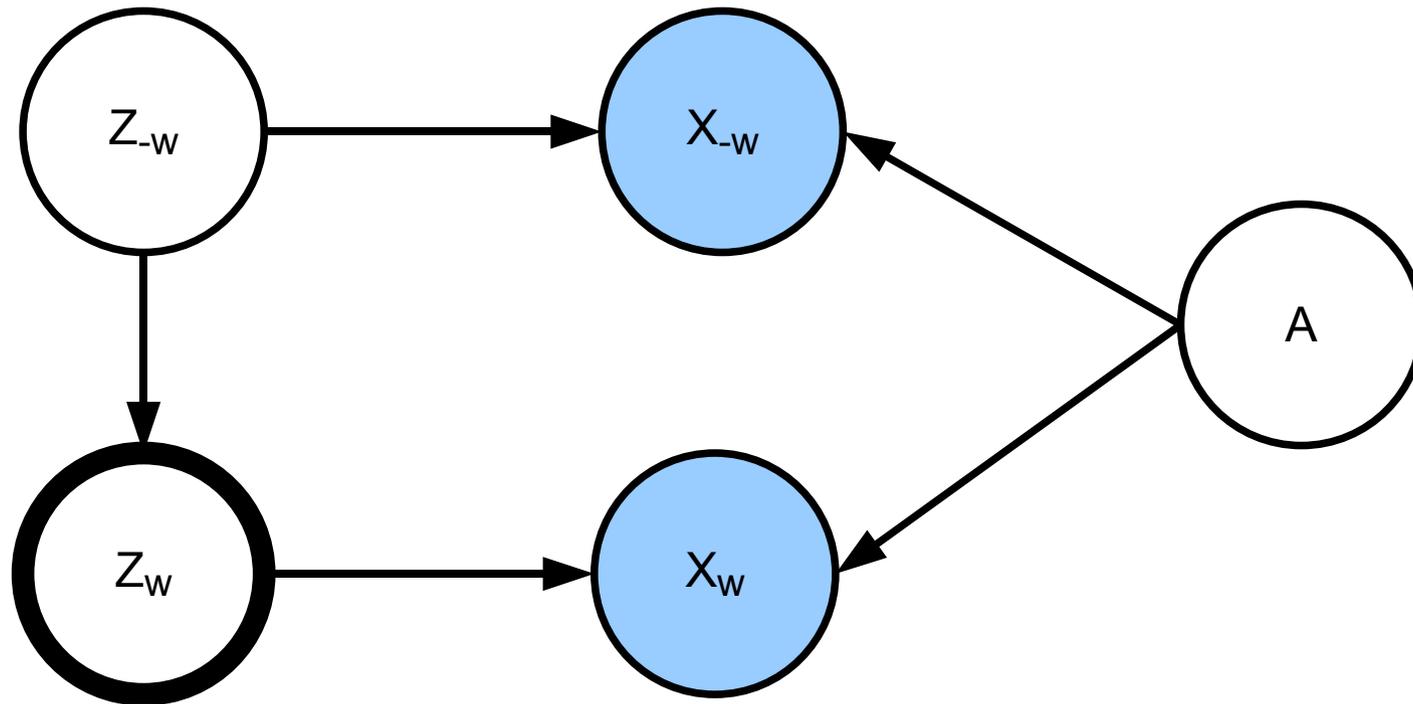
# Full Model

Data
Matrix
X

Features
Present
Z

Feature
Values
A

# Full Model



Data Matrix X

Features Present Z

Feature Values A

$X_{-w}$ { ... $= Z_{-w}$ { ... $*$ ...

$X_w$ { ... $Z_w$ { ...

Note: this is not Blocked Gibbs Sampling!

UNIVERSITY OF
CAMBRIDGE

# The Graphical Model

# Basic Sampling

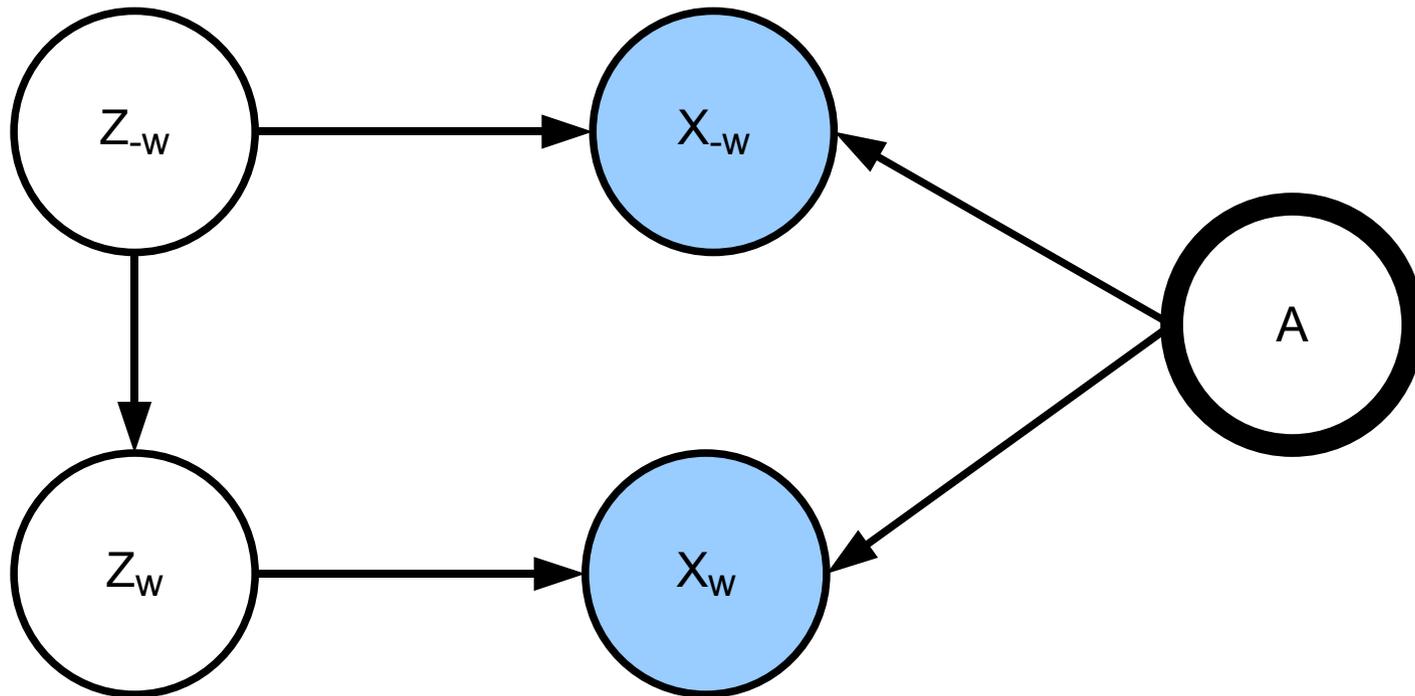First sample $Z_W | X, A, Z_{-W}$

UNIVERSITY OF
CAMBRIDGE

# Basic Sampling

First sample $Z_W | X, A, Z_{-W}$ and then $Z_{-W} | X, A, Z_W$
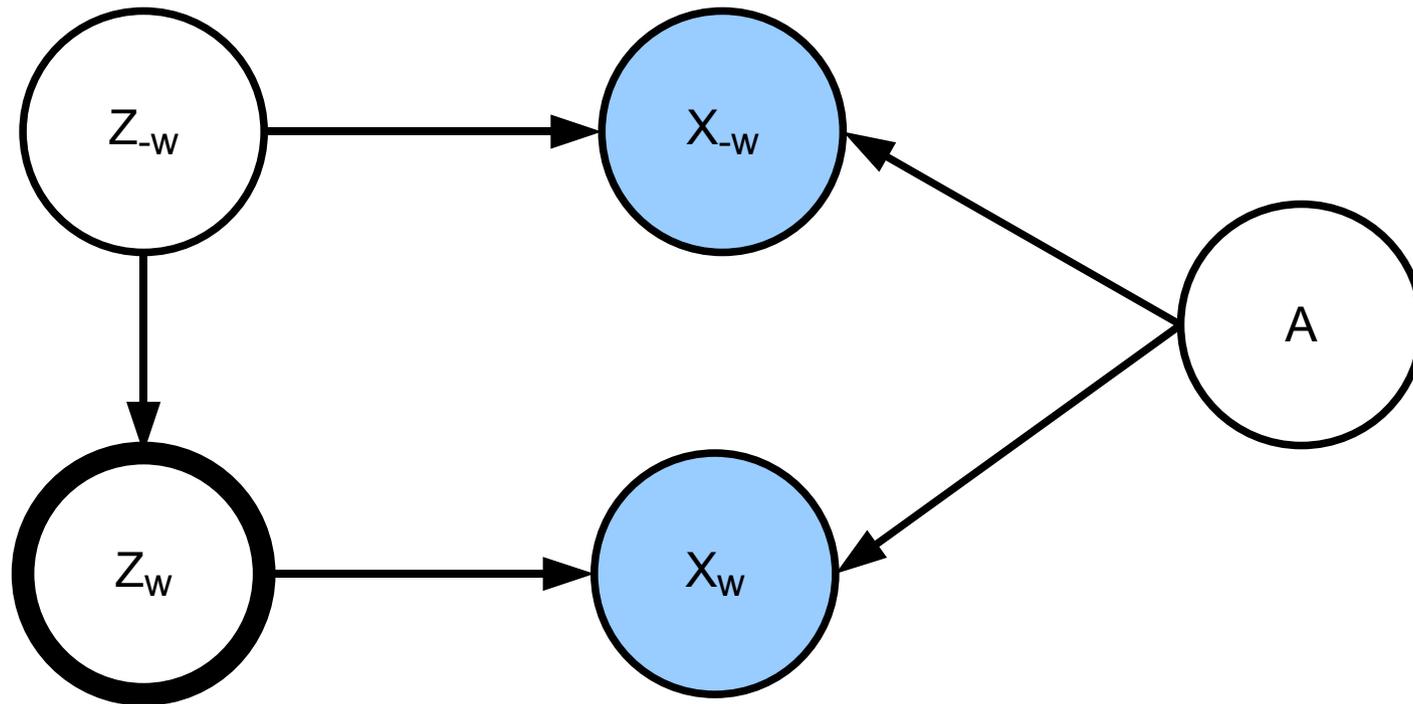
# Basic Sampling

First sample $Z_w | X, A, Z_{-w}$ and then $Z_{-w} | X, A, Z_w$ and then $A | Z, X$ ...
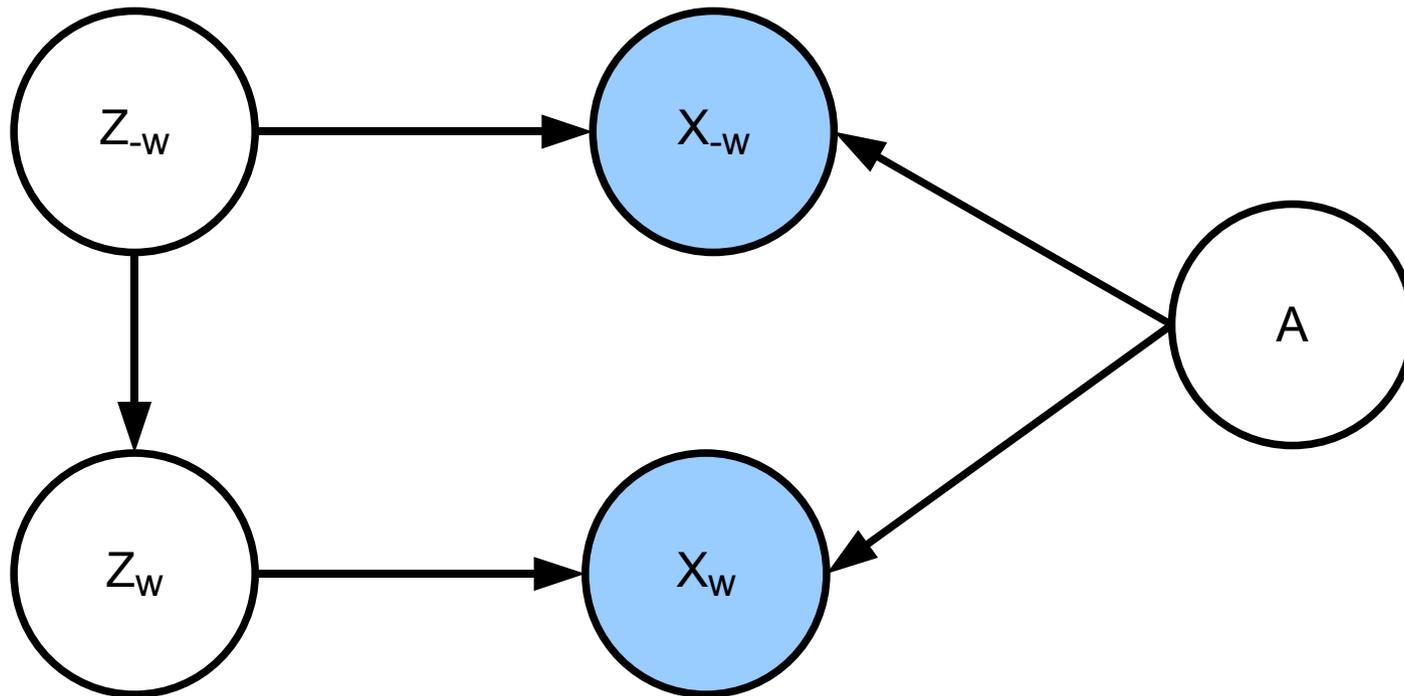
# Basic Sampling

First sample $Z_W | X, A, Z_{-W}$ and then $Z_{-W} | X, A, Z_W$ and then $A | Z, X$ and then $Z_W | X, A, Z_{-W}$ ...
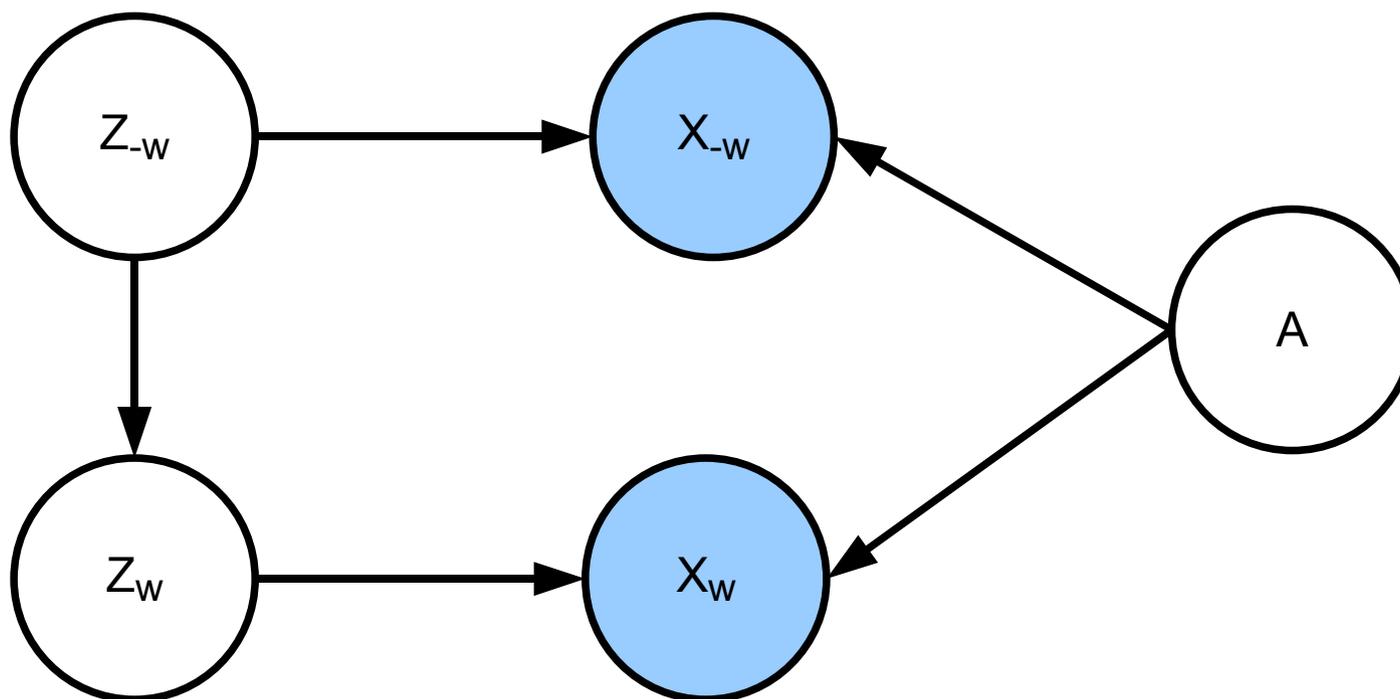
# Basic Sampling

Advantage: Each iteration is fast to compute.

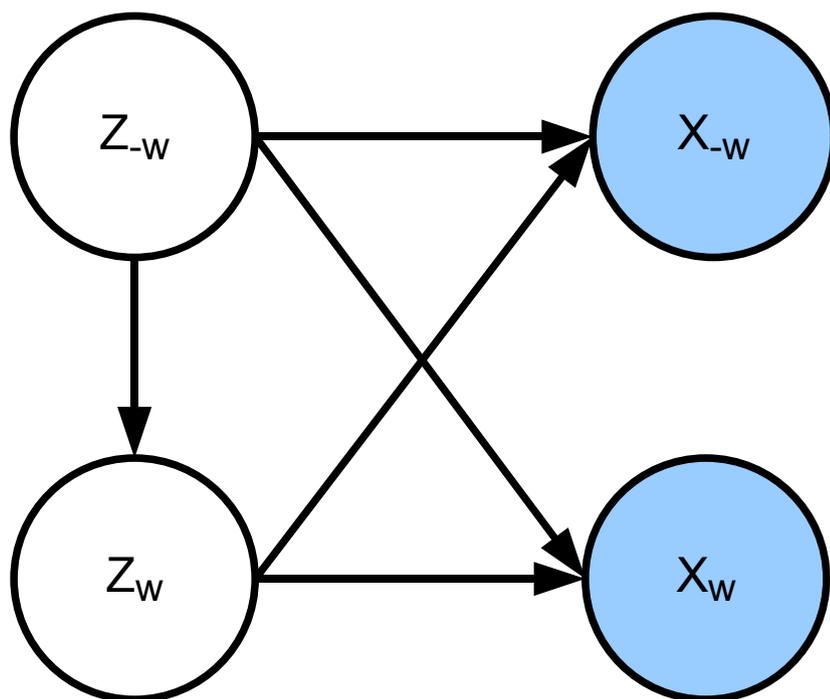Disadvantage: Often slow to mix.

# Collapsed Gibbs Sampling

Since we can compute P(X|Z), integrate out A

UNIVERSITY OF
CAMBRIDGE

# Collapsed Gibbs Sampling

Since we can compute P(X|Z), integrate out A

UNIVERSITY OF
CAMBRIDGE

# Collapsed Gibbs Sampling

Sample each Z in turn, as before

UNIVERSITY OF
CAMBRIDGE

# Collapsed Gibbs Sampling

Sample each Z in turn, as before

UNIVERSITY OF
CAMBRIDGE

# Collapsed Gibbs Sampling

Advantage: Faster to mix.

Disadvantage: Inference no longer scales!

UNIVERSITY OF
CAMBRIDGE

# Our solution: Accelerated Sampling

Keep a posterior on A.  Observations stay independent!
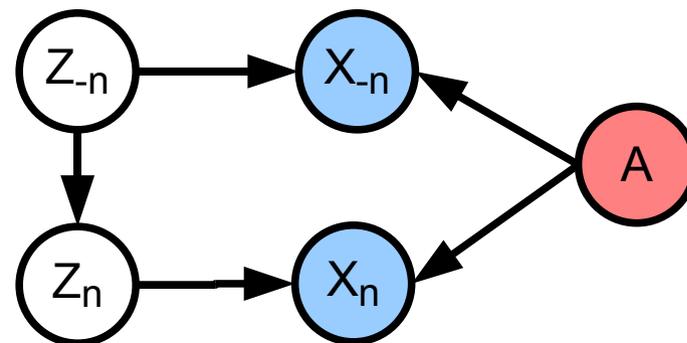
# More formally: Consider one element

$$P(Z_{nk}=1|Z_{-nk}, X) \propto$$



$$P(Z_{nk}=1|Z_{-nk}) P(X|Z)$$

$$P(Z_{nk}=1|Z_{-nk}) \int_A P(X|Z, A) P(A) \, dA$$

$$P(Z_{nk}=1|Z_{-nk}) \int_A P(X_n|Z_n, A) P(X_{-n}|Z_{-n}, A) P(A) \, dA$$

$$P(Z_{nk}=1|Z_{-nk}) \int_A P(X_n|Z_n, A) \, {\color{red} P(A|Z_{-n}, X_{-n})} \, dA$$

UNIVERSITY OF
CAMBRIDGE

# More formally: Consider one element

$$P(Z_{nk}=1|Z_{-nk}, X) \propto$$

Bayes Rule



$$P(Z_{nk}=1|Z_{-nk}) P(X|Z)$$

$$P(Z_{nk}=1|Z_{-nk}) \int_A P(X|Z, A) P(A) \, dA$$

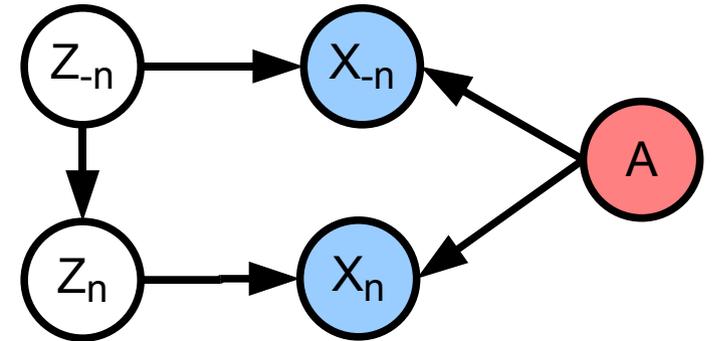$$P(Z_{nk}=1|Z_{-nk}) \int_A P(X_n|Z_n, A) P(X_{-n}|Z_{-n}, A) P(A) \, dA$$

$$P(Z_{nk}=1|Z_{-nk}) \int_A P(X_n|Z_n, A) P(A|Z_{-n}, X_{-n}) \, dA$$

UNIVERSITY OF CAMBRIDGE

# More formally: Consider one element

$$P(Z_{nk}=1|Z_{-nk}, X) \propto$$



$$P(Z_{nk}=1|Z_{-nk}) P(X|Z)$$

$$P(Z_{nk}=1|Z_{-nk}) \int_A P(X|Z, A) P(A) dA$$

Joints and conditionals

$$P(Z_{nk}=1|Z_{-nk}) \int_A P(X_n|Z_n, A) P(X_{-n}|Z_{-n}, A) P(A) dA$$

$$P(Z_{nk}=1|Z_{-nk}) \int_A P(X_n|Z_n, A) P(A|Z_{-n}, X_{-n}) dA$$

UNIVERSITY OF CAMBRIDGE
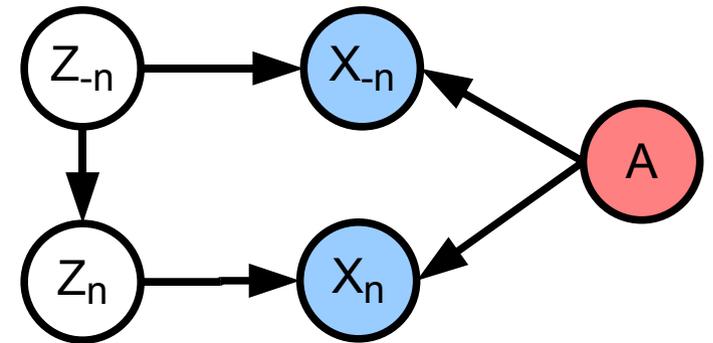
# More formally: Consider one element

$$P(Z_{nk}{=}1|Z_{-nk}, X) \propto$$

$$P(Z_{nk}{=}1|Z_{-nk}) P(X|Z)$$

$$P(Z_{nk}{=}1|Z_{-nk}) \int_A P(X|Z, A) P(A) \, dA$$

$$P(Z_{nk}{=}1|Z_{-nk}) \int_A P(X_n|Z_n, A) P(X_{-n}|Z_{-n}, A) P(A) \, dA$$

$$P(Z_{nk}{=}1|Z_{-nk}) \int_A P(X_n|Z_n, A) P(A|Z_{-n}, X_{-n}) \, dA$$

Bayes Rule

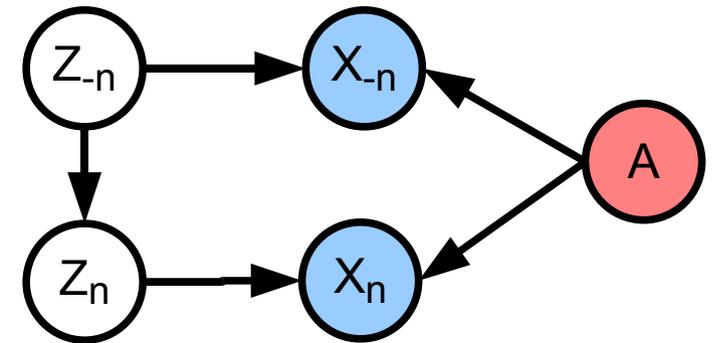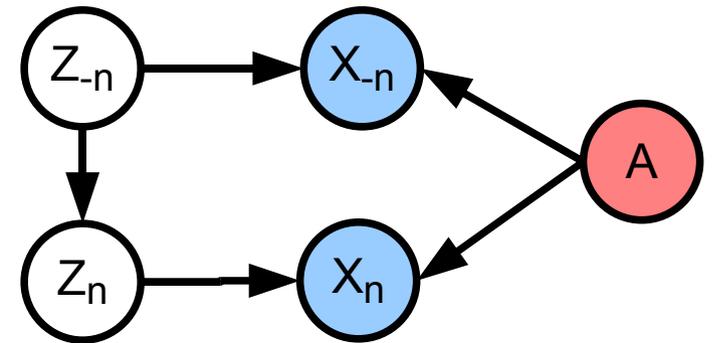UNIVERSITY OF CAMBRIDGE

# More formally: Consider one element

$$P(Z_{nk}=1|Z_{-nk}, X) \propto$$



$$P(Z_{nk}=1|Z_{-nk}) P(X|Z)$$

$$P(Z_{nk}=1|Z_{-nk}) \int_A P(X|Z, A) P(A) \, dA$$

$$P(Z_{nk}=1|Z_{-nk}) \int_A P(X_n|Z_n, A) P(X_{-n}|Z_{-n}, A) P(A) \, dA$$
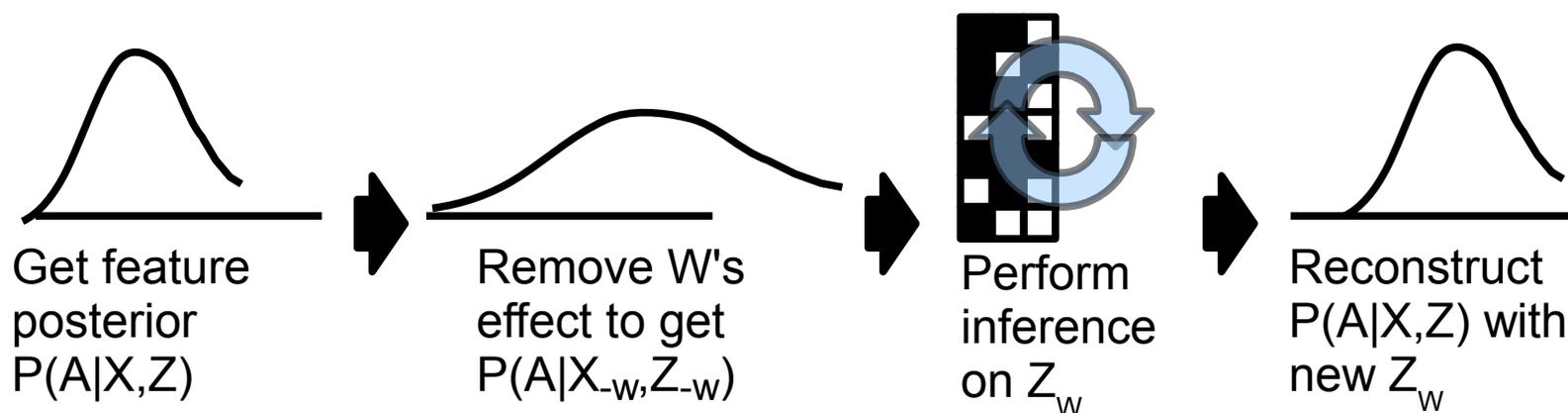
$$P(Z_{nk}=1|Z_{-nk}) \int_A P(X_n|Z_n, A) \, {\color{red} P(A|Z_{-n}, X_{-n})} \, dA$$

EXACT!

UNIVERSITY OF
CAMBRIDGE

# Accelerated Gibbs Sampling

1. Initialise some Z, feature posterior

2. For each window of observations W



Get feature posterior $P(A|X,Z)$ → Remove W's effect to get $P(A|X_{-w},Z_{-w})$ → Perform inference on $Z_w$ → Reconstruct $P(A|X,Z)$ with new $Z_w$

Considerations: how many observations should we consider at once? Depends on the cost of computing $P(A|X,Z)$ and $P(X|Z,A)$, numerical errors.
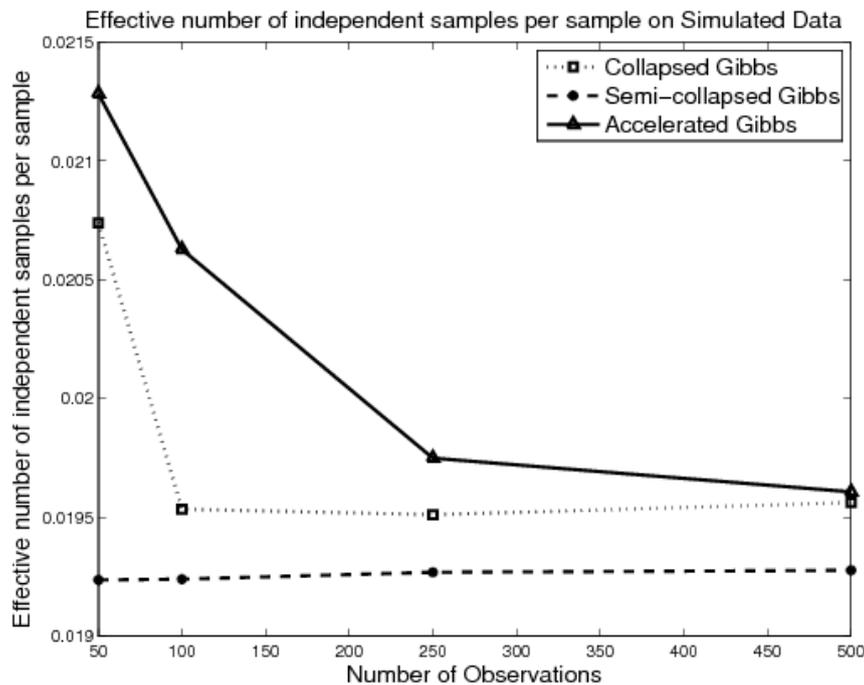
# Details for the IBP Model

If the prior on A, noise is Gaussian, then

- Posterior on A is Gaussian.

- Posterior can be updated with rank-one updates.
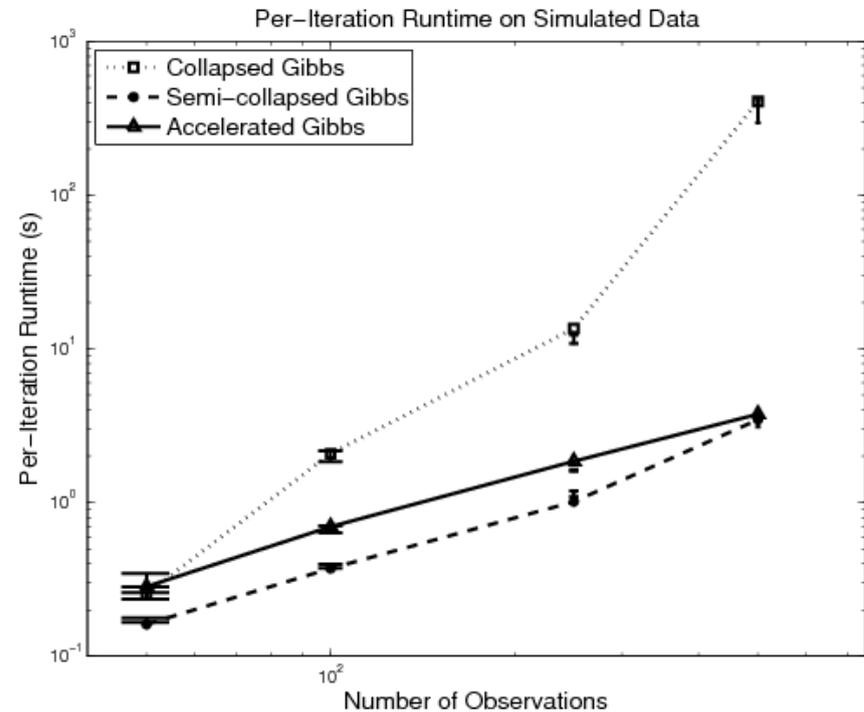
- Optimal window is 1.

Also, intelligently choosing to represent Gaussians in information form $(h, \Sigma^{-1})$ or covariance form $(\mu, \Sigma)$ helps maintain numerical precision.  Details in the paper.

UNIVERSITY OF
CAMBRIDGE

# Experiments on Synthetic Data

Data generated from the prior; D=10, N = {50,100,250, 500}.



Mixing similar to
collapsed sampler

Runtime similar to
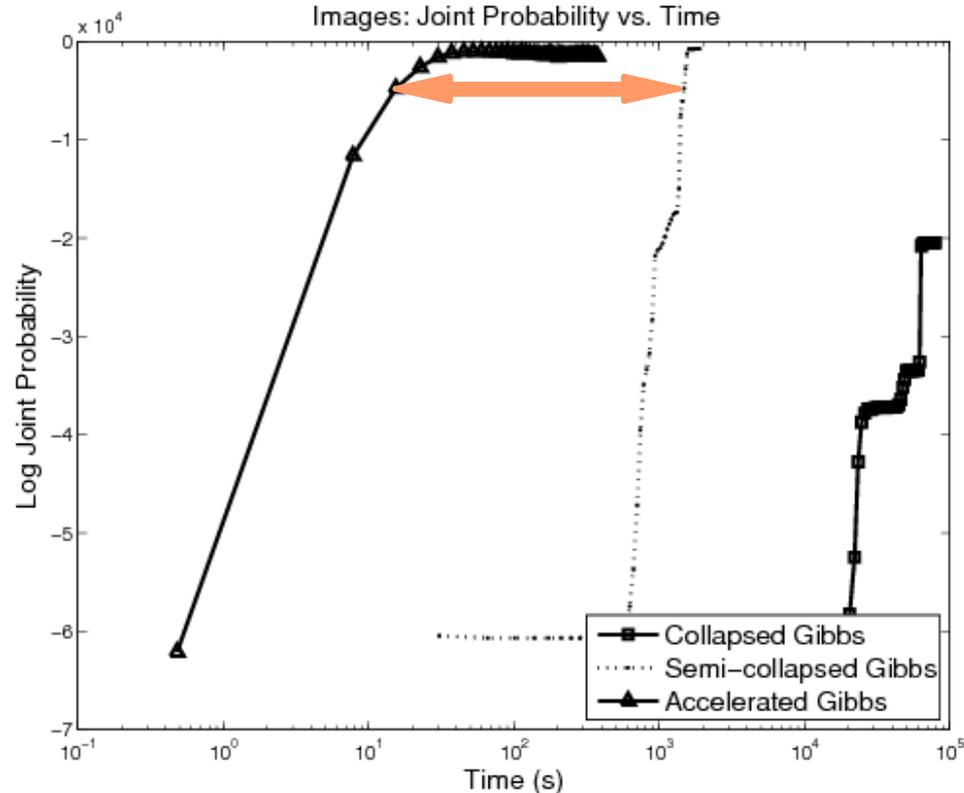semi-collapsed sampler

UNIVERSITY OF
CAMBRIDGE

# Experiments on Smaller Datasets

## D=36, N = 1000



## D=1024, N = 722



Reach mode orders of
magnitude faster!

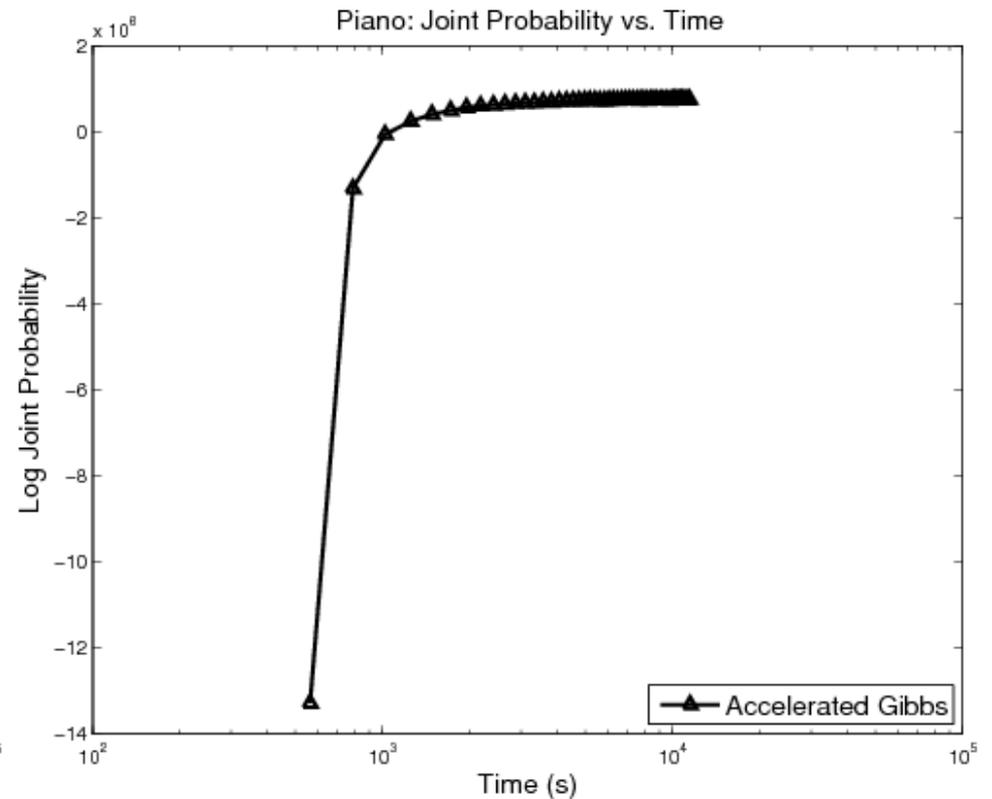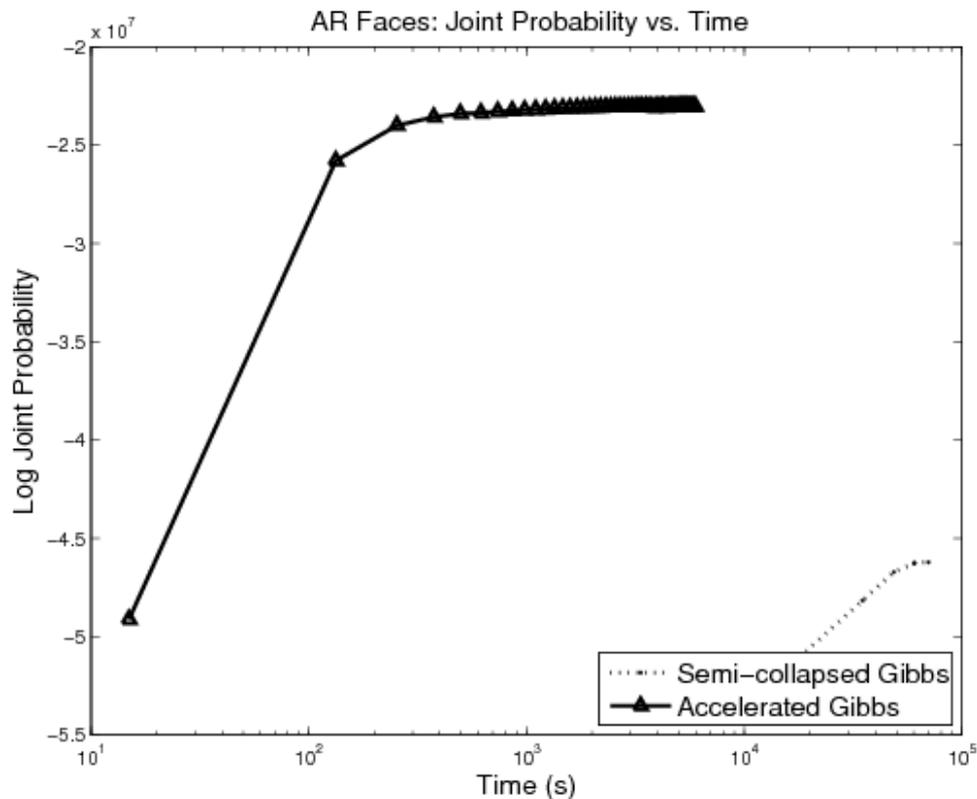UNIVERSITY OF
CAMBRIDGE

# Experiments on Larger Datasets

D=1598, N = 2600                    D=161, N = 10000



Standard samplers
become impractical...

# Returning to an age-old question...

To marginalize or not marginalize, that is the question:

Whether 'tis more tractable for the sampler to suffer the hills and valleys of local optima,

Or to take expectations against a set of variables, and by integrating collapse them?

# Returning to an age-old question...

To marginalize or not marginalize, that is the question:

Whether 'tis more tractable for the sampler to suffer the hills and valleys of local optima,

Or to take expectations against a set of variables, and by integrating collapse them?
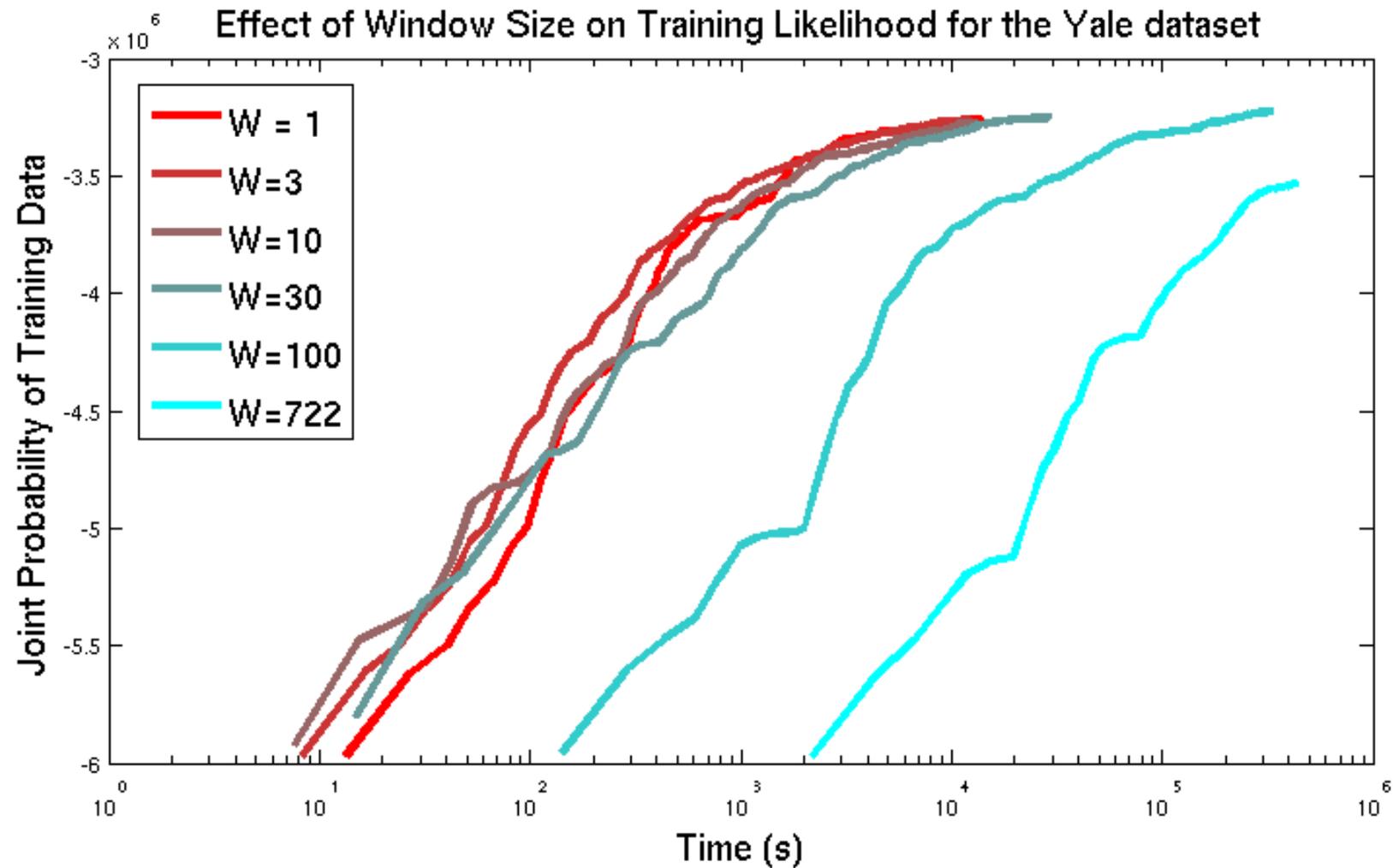
*In answer: of a third example...*

# Conclusions

- Maintaining a posterior within a sampler allows us to perform fast inference in an important class of models

- In particular, our approach allows us to scale inference to large Indian Buffet Process models.

... code available on my website:
http://mlg.eng.cam.ac.uk/finale/wiki

UNIVERSITY OF
CAMBRIDGE

# Effect of Window Size

# Experiments on Real Data



OriginalReconstructed Parts

UNIVERSITY OF
CAMBRIDGE

# EEG Dataset



EEG: Joint Probability vs. Time

UNIVERSITY OF
CAMBRIDGE