

# Probabilistic Deterministic Infinite Automata

David Pfau, Nicholas Bartlett, Frank Wood

{pfau@neurotheory, bartlett@stat, fwood@stat}.columbia.edu



COLUMBIA UNIVERSITY  
IN THE CITY OF NEW YORK

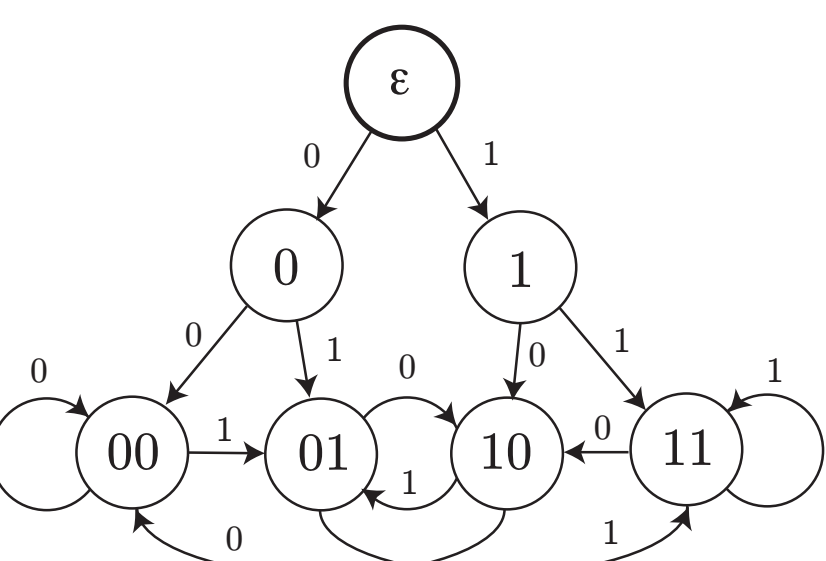
## Overview

A Probabilistic Deterministic Infinite Automata (PDIA) is a Probabilistic Deterministic Finite Automata (PDFA)[1] with an unbounded number of states. We take a nonparametric Bayesian approach to PDIA inference.

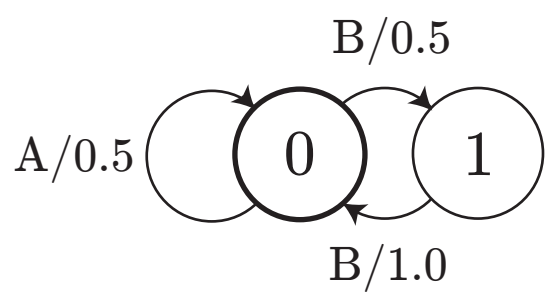
PDIA are an attractive compromise between the computational costs of Hidden Markov Models and the storage requirements of smoothed Markov models for predicting sequence data.

## Finite Automata

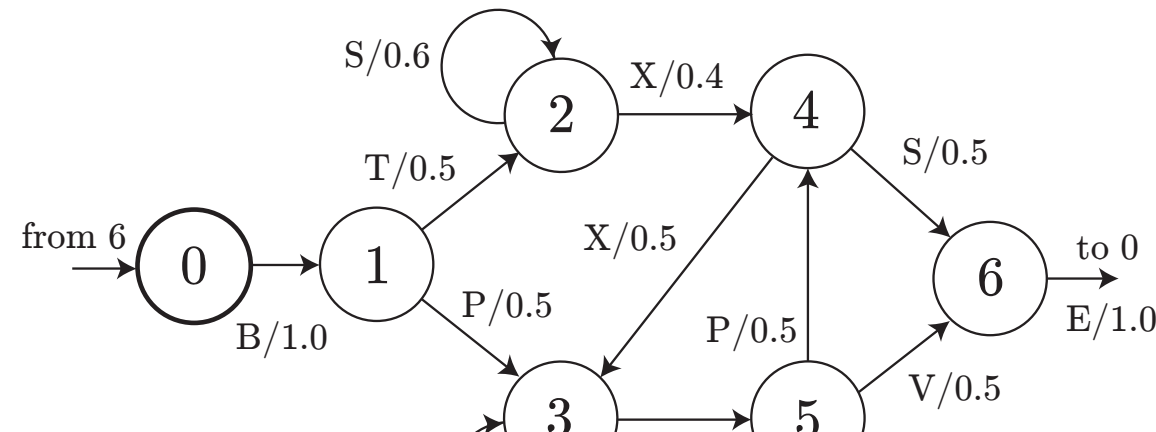
### Probabilistic Deterministic Finite Automata



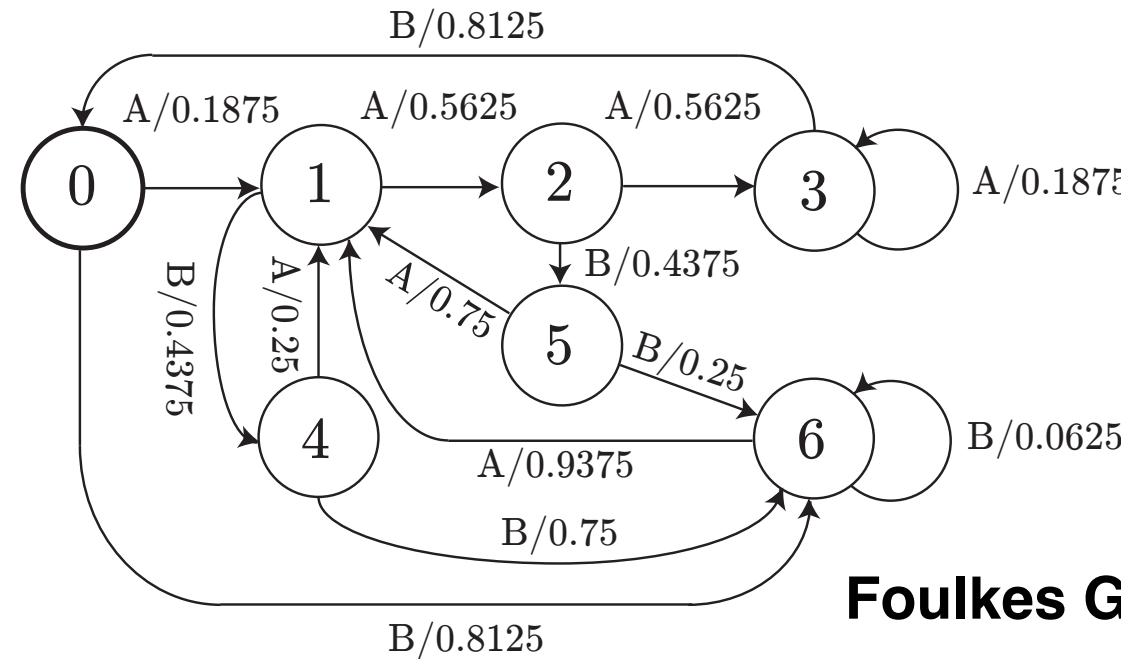
Trigram as DFA (without probability)



Even Process



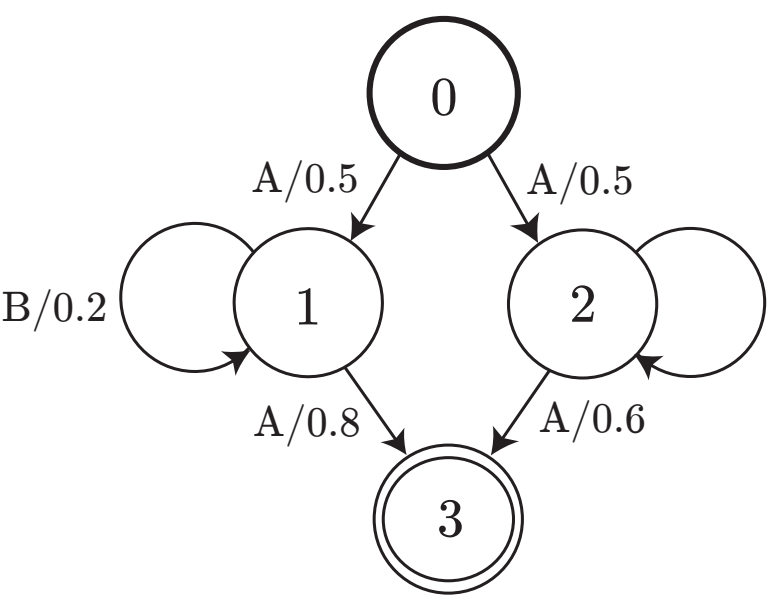
Reber Grammar



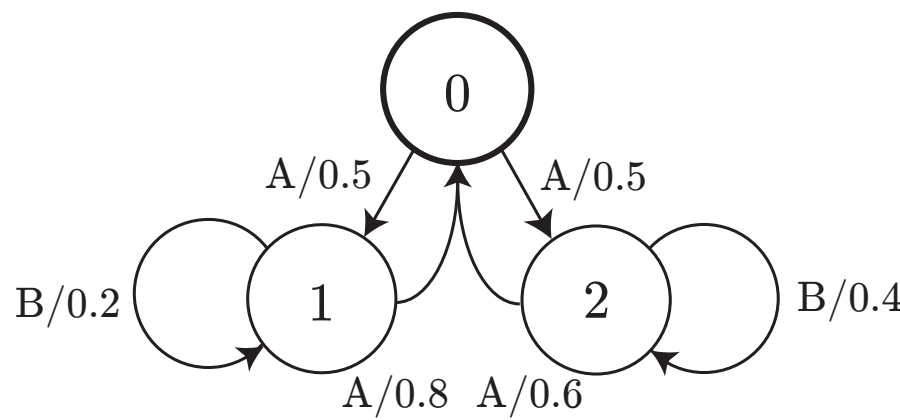
Foulkes Grammar

### Probabilistic Nondeterministic Finite Automata

(a)



(b)



(a) in mixtures of PDFA but not PDFA

(b) in PNFA but not mixtures of PDFA

Given a finite training sequence, the posterior distribution over PDIA parameters is an infinite mixture of PDFAs.

Samples for this distribution drawn via MCMC form a finite mixture approximating this posterior.

The class of models we deal with, from least to most general -

$n^{\text{th}}$ -order Markov  $\subset$  PDFA  $\subset$  finite mixture of PDFA  $\subset$  PNFA = HMM

(Technically, HMM = PNFA without final state)

## Generative Model

### A Prior over PDFA with a bounded number of states

$$\begin{array}{c} \delta \\ \sigma_0 \quad \sigma_1 \quad \sigma_2 \\ q_0 \begin{bmatrix} q_3 & q_1 & q_{12} \\ \bullet & q_3 & \bullet \\ q_6 & q_4 & \bullet \\ q_3 & \bullet & q_5 \\ q_3 & q_{42} & q_5 \\ \bullet & q_1 & q_2 \end{bmatrix} \\ q_1 \\ q_2 \\ q_3 \\ q_4 \\ q_5 \\ \vdots \end{array}$$

$$\begin{array}{c} \pi \\ \sigma_0 \quad \sigma_1 \quad \sigma_2 \\ q_0 \begin{bmatrix} \text{(iid probability vector)} \\ \text{(iid probability vector)} \\ \vdots \end{bmatrix} \\ q_1 \\ \vdots \end{array}$$

$$\begin{array}{c} q_0 \quad q_1 \quad q_3 \quad q_3 \quad q_3 \quad q_5 \quad q_2 \quad q_4 \quad q_5 \quad q_2 \dots \\ 1 \quad 1 \quad 0 \quad 0 \quad 2 \quad 2 \quad 1 \quad 2 \quad 1 \quad 0 \dots \end{array}$$

This prior biases towards state reuse in two ways: the global bias toward some states due to  $\mu$  and the column-specific bias towards certain states when the same symbol is emitted due to  $\phi_j$ .

### A Prior over infinite models

The limit as  $|Q| \rightarrow \infty$  is a PDIA prior. The two-level Dirichlet construction of the finite model becomes a Hierarchical Dirichlet Process (HDP)[2] which we give geometric base distribution. When  $\mu$  and  $\phi_j$  are marginalized out,  $\delta_{ij}$  are exchangeable.

## Inference

MCMC sampler for  $\delta_{ij} | \delta_{-ij}, x_{0:T}, \alpha, \alpha_0, \beta$

Marginalizing out  $\pi$  in likelihood gives form that only depends on counts

$$p(x_{0:T} | \delta, c, \beta) = \prod_{i=0}^{|Q|-1} \frac{\Gamma(\beta)}{\Gamma(\frac{\beta}{|\Sigma|})^{|\Sigma|}} \frac{\prod_{j=1}^{|\Sigma|} \Gamma(\frac{\beta}{|\Sigma|} + c_{ij})}{\Gamma(\beta + \sum_{j=1}^{|\Sigma|} c_{ij})} \quad c_{ij} = \sum_{t=0}^T \mathbf{1}_{ij}(\xi_t, x_t)$$

Propose  $\delta_{ij}^*$  from  $\delta_{ij}^* | \delta_{-ij}, \alpha, \alpha_0$ , easy if  $|Q| < \infty$ , use *Chinese Restaurant Franchise Process*[2] if  $|Q| = \infty$

Accept with probability  $\min \left( 1, \frac{p(x_{0:T} | \delta_{ij}^*, \delta_{-ij})}{p(x_{0:T} | \delta_{ij}, \delta_{-ij})} \right)$

If  $c_{ij} = 0$ ,  $\delta_{ij}$  can be ignored

## Natural Language, DNA Prediction

|     | PDIA  | PDIA-MAP | HMM-EM | bigram | trigram | 4-gram | 5-gram | 6-gram | SSM     |
|-----|-------|----------|--------|--------|---------|--------|--------|--------|---------|
| AIW | 5.13  | 5.46     | 7.89   | 9.71   | 6.45    | 5.13   | 4.80   | 4.69   | 4.78    |
|     | 365.6 | 379      | 52     | 28     | 382     | 2,023  | 5,592  | 10,838 | 19,358  |
| DNA | 3.72  | 3.72     | 3.76   | 3.77   | 3.75    | 3.74   | 3.73   | 3.72   | 3.56    |
|     | 64.7  | 54       | 19     | 5      | 21      | 85     | 341    | 1,365  | 314,166 |

Top rows: perplexity of held out data. Bottom rows: number of states

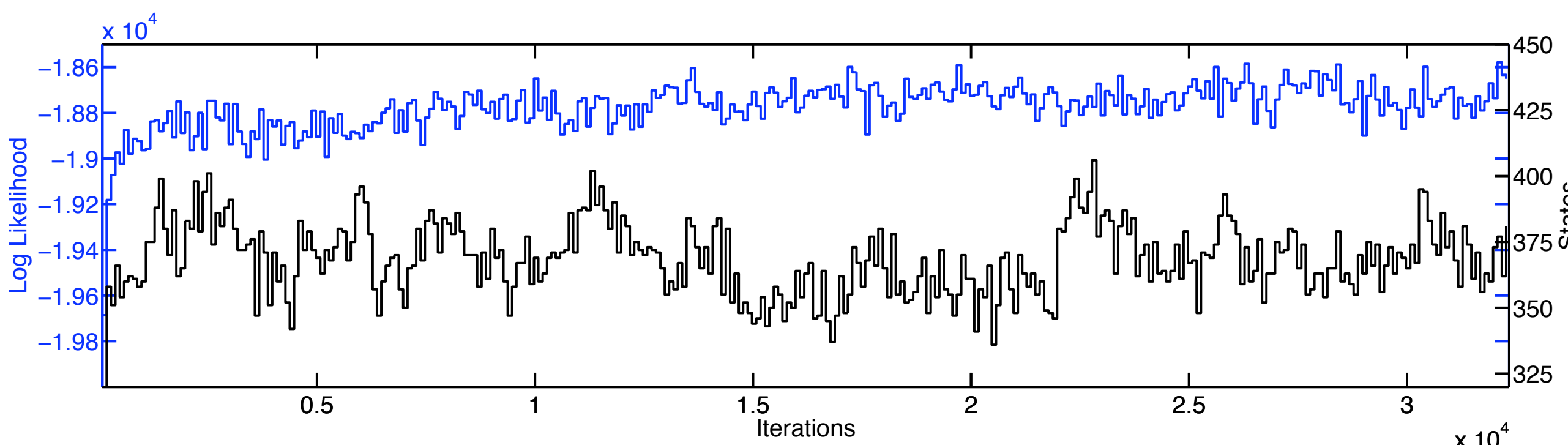
Predictions from average of many samples better than MAP sample

Low model complexity relative to other models with the same predictive performance

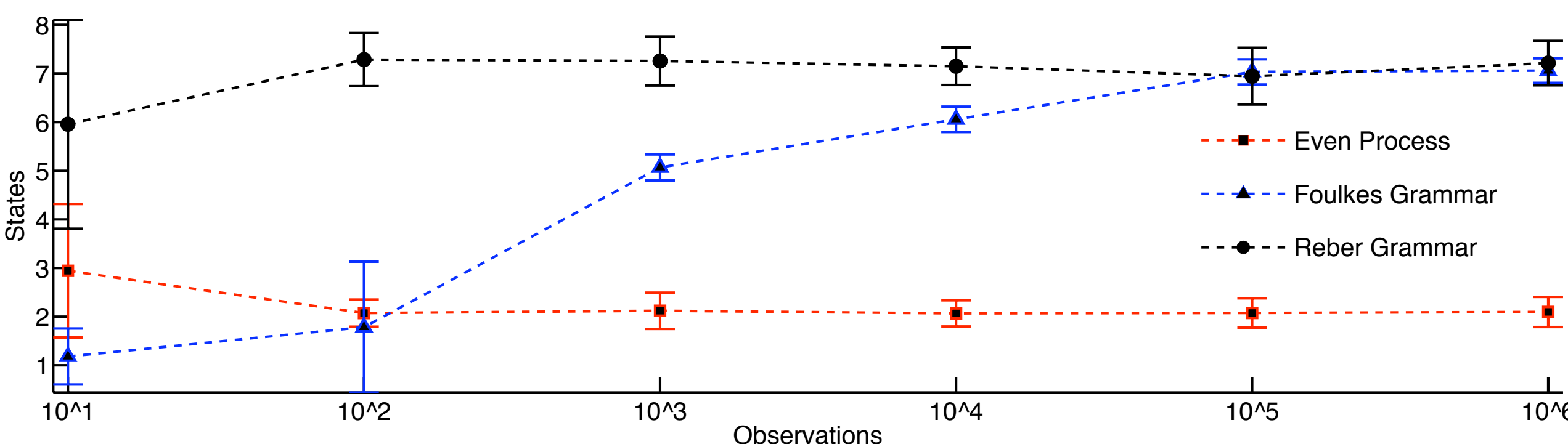
Alice in Wonderland: 10k training characters, 4k test “*alice was beginning to...*”

Mouse DNA: 150k train, 50k test “*CGTATATGCGCC...*”

Controls: EM-trained HMM, HPYP smoothed n-gram[3], sequentially-trained sequence memoizer[4]



## Synthetic Grammar Induction



## Future Directions

Evaluation on larger data sets

Sampler that proposes “similar” PDFA - spilt-merge sampling?

How to smooth emission distributions?

## References

- [1] M.O. Rabin. Probabilistic Automata. *Information and Control*, 6(3):230-245, 1963
- [2] Y. W. Teh, M. I. Jordan, M. J. Beal, D. M. Blei. Hierarchical Dirichlet Processes. *Journal of the American Statistical Association*. 101(476):1566-1581, 2006
- [3] Y. W. Teh. A hierarchical Bayesian language model based on Pitman-Yor processes. *Proceedings of the ACL*. 985-992, 2006
- [4] F. Wood, C. Archambeau, J. Gasthaus, L. James, Y. W. Teh. A Stochastic Memoizer for Sequence Data. *Proceedings of the 26th ICML*. 1129-1136, 2009

This research was supported by the NSF GRFP