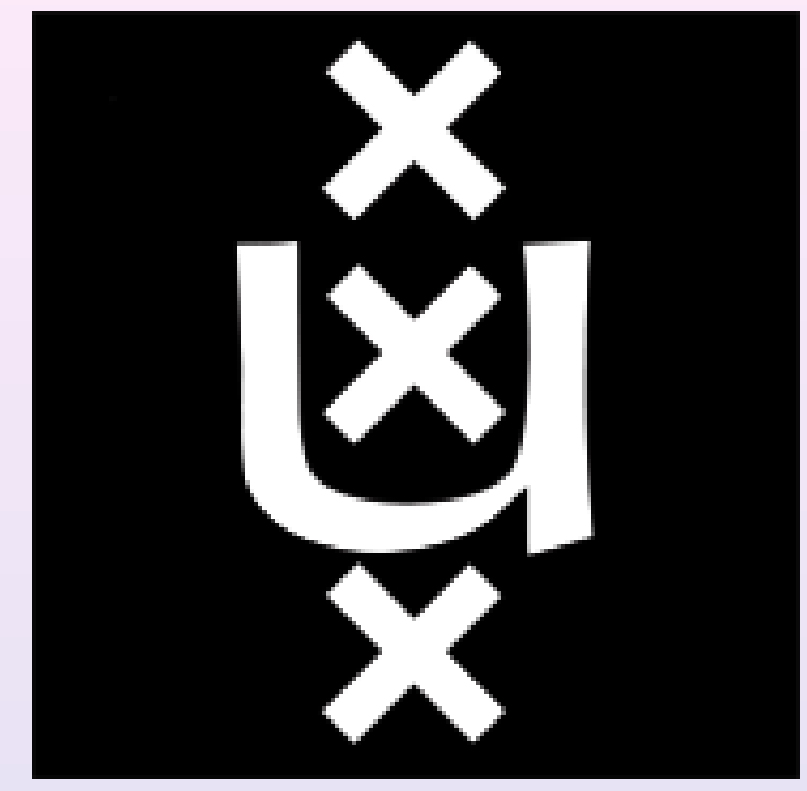




# Stochastic Collapsed Variational Bayesian Inference for Latent Dirichlet Allocation



James Foulds<sup>1</sup>, Levi Boyles<sup>1</sup>, Christopher DuBois<sup>2</sup>, Padhraic Smyth<sup>1</sup>, Max Welling<sup>3</sup>

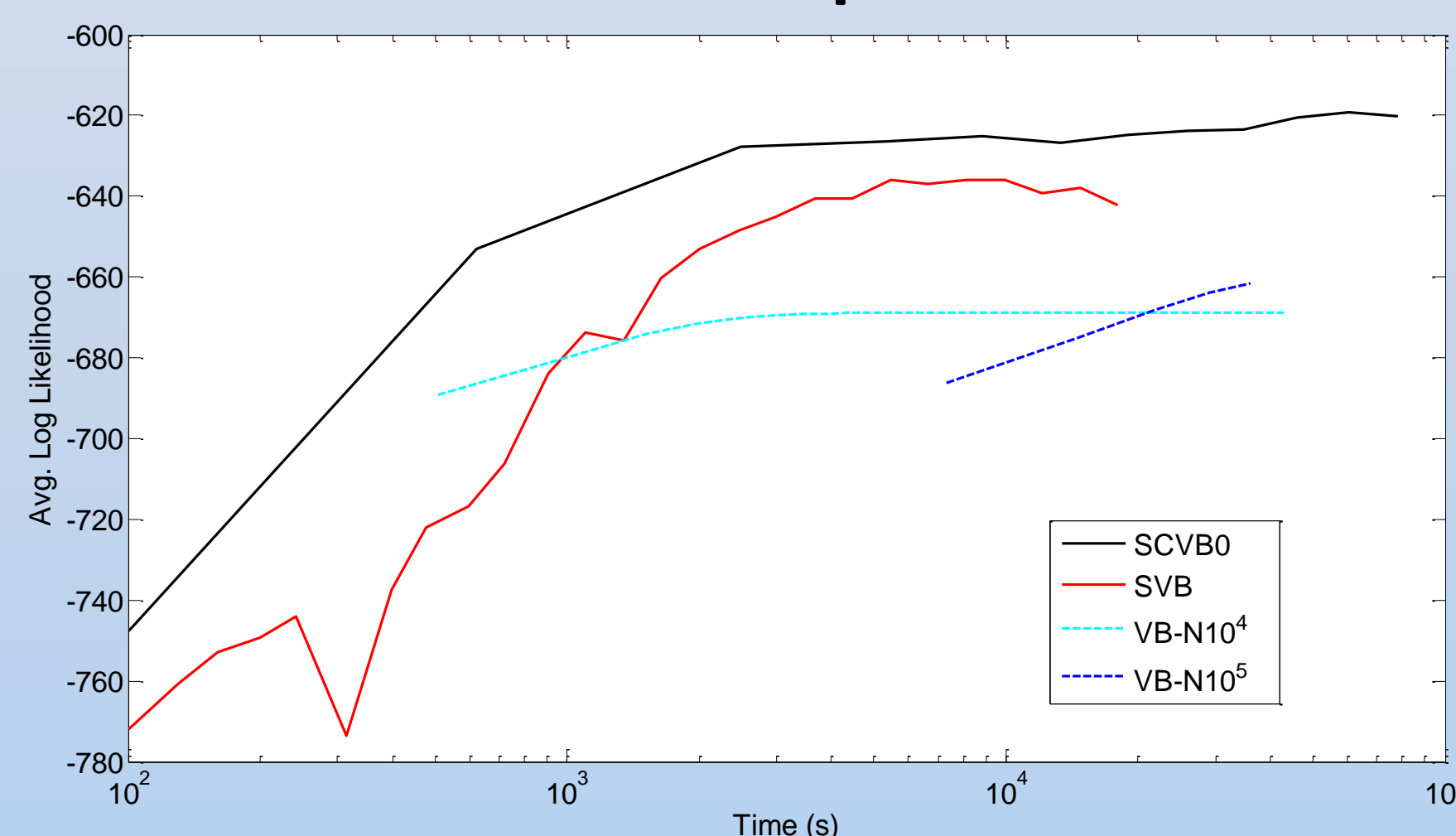
<sup>1</sup>University of California Irvine, Computer Science

<sup>2</sup>University of California Irvine, Statistics

<sup>3</sup>University of Amsterdam, Computer Science

## Motivation

### LDA on Wikipedia



### Stochastic Optimization for ML

- **Batch algorithms**
  - While (not converged)
    - Process the **entire** dataset
    - Update parameters
- **Stochastic algorithms**
  - While (not converged)
    - Process a **subset** of the dataset
    - Estimate quantities needed for an update, and extrapolate
    - Update parameters

### Collapsed Variational Bayes for LDA

- Maintain variational distributions for the topic of each token
- Mean field assumption
- CVB0 (Asuncion et al., 2009)

$$\gamma_{ijk} \propto \frac{N_{w_{ij}k}^{\Phi} + \eta_{w_{ij}}}{N_k^Z + \sum_w \eta_w} (N_{jk}^{\Theta} + \alpha_k)$$

$$N_k^Z \triangleq \sum_{ij} \gamma_{ijk}$$

$$N_{jk}^{\Theta} \triangleq \sum_i \gamma_{ijk}$$

$$N_{wk}^{\Phi} \triangleq \sum_{ij: w_{ij}=w} \gamma_{ijk}$$

### Advantages of the collapsed representation

- **Simpler**, faster and fewer update equations
- **Better mixing** for Gibbs sampling
- **Better variational bound** for VB

## Stochastic CVB0

### How to estimate CVB0 statistics?

- Pick a random word (i,j) from the corpus

$$N_k^Z \triangleq \sum_{ij} \gamma_{ijk}$$

$$E[N_k^Z] = \text{Total words in the corpus} \times \gamma_{ijk}$$

- In an online algorithm, we cannot store the variational parameters, but we can *update them*

$$\gamma_{ijk} \propto \frac{N_{w_{ij}k}^{\Phi} + \eta_{w_{ij}}}{N_k^Z + \sum_w \eta_w} (N_{jk}^{\Theta} + \alpha_k)$$

- Keep an online average of the CVB0 statistics

$$N_j^{\Theta} := (1 - \rho_t^{\Theta}) N_j^{\Theta} + \rho_t^{\Theta} \hat{N}_j^{\Theta}$$

$$N^{\Phi} := (1 - \rho_t^{\Phi}) N^{\Phi} + \rho_t^{\Phi} \hat{N}^{\Phi}$$

$$N^Z := (1 - \rho_t^Z) N^Z + \rho_t^Z \hat{N}^Z$$

### The Full Algorithm

- Optional burn-in passes per document
- Minibatches
- Operating on sparse counts

$$N_i^{\Theta} := (1 - \rho_t^{\Theta}) N_i^{\Theta} + \rho_t^{\Theta} C_j \gamma_{ij}$$

$$N^{\Phi} := (1 - \rho_t^{\Phi}) N^{\Phi} + \rho_t^{\Phi} \hat{N}^{\Phi}$$

$$N^Z := (1 - \rho_t^Z) N^Z + \rho_t^Z \hat{N}^Z$$

$$\text{where } \hat{N}^{\Phi} = \frac{C}{|M|} \sum_{ij \in M} \mathbf{Y}^{(ij)} \text{ and } \hat{N}^Z = \frac{C}{|M|} \sum_{ij \in M} \gamma_{ij}$$

$$N_j^{\Theta} := (1 - \rho_t^{\Theta})^{m_{aj}} N_j^{\Theta} + C_j \gamma_{aj} (1 - (1 - \rho_t^{\Theta})^{m_{aj}})$$

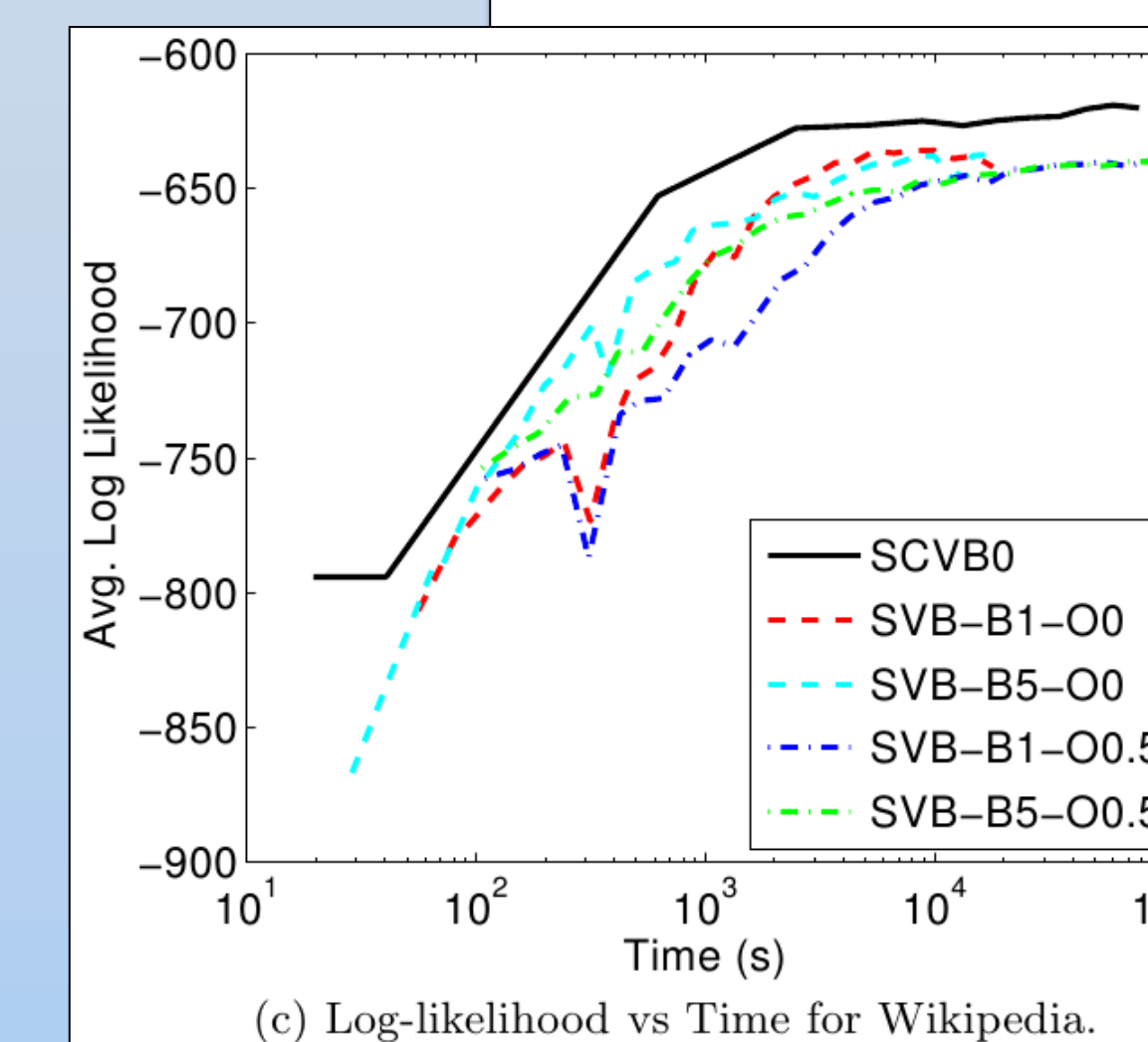
- Randomly initialize  $N^{\Phi}$ ,  $N^{\Theta}$ ;  $N^Z := \sum_w N_w^{\Phi}$
- For each minibatch  $M$ 
  - $\hat{N}^{\Phi} := \mathbf{0}$ ;  $\hat{N}^Z := \mathbf{0}$
  - For each document  $j$  in  $M$ 
    - For zero or more “burn-in” passes
      - For each token  $i$ 
        - Update  $\gamma_{ij}$
        - Update  $N_i^{\Theta}$
    - For each token  $i$ 
      - Update  $\gamma_{ij}$
      - Update  $N_i^{\Theta}$
      - $\hat{N}_{w_{ij}}^{\Phi} := \hat{N}_{w_{ij}}^{\Phi} + \frac{C}{|M|} \gamma_{ij}$
      - $\hat{N}^Z := \hat{N}^Z + \frac{C}{|M|} \gamma_{ij}$
  - Update  $N^{\Phi}$
  - Update  $N^Z$

SCVB0 is a Robbins Monro **stochastic approximation** algorithm for finding the fixed points of (a variant of) CVB0

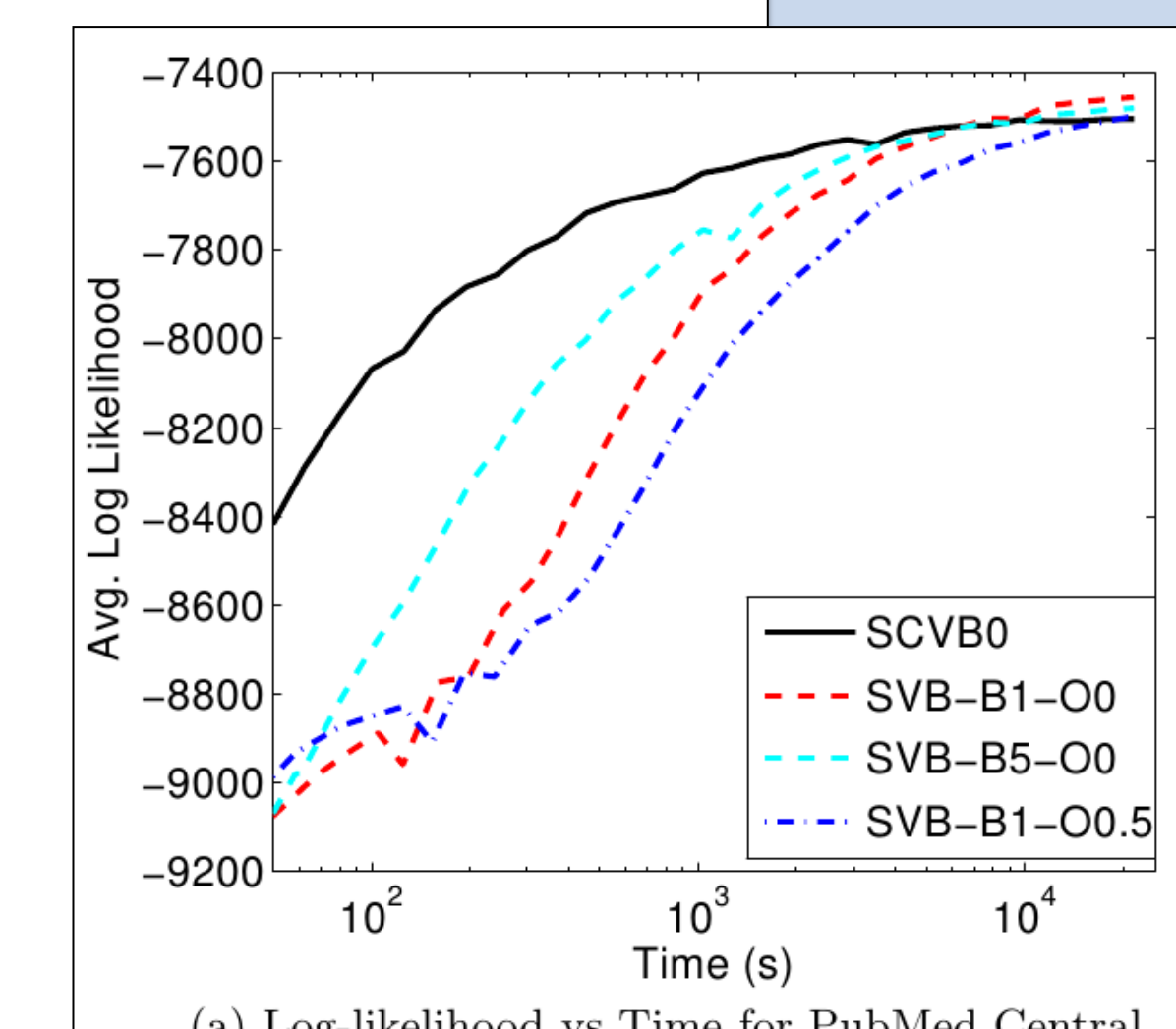
Theorem: with an appropriate sequence of step sizes, **SCVB0 converges to a stationary point of the MAP**, with adjusted hyper-parameters

## Experimental Results

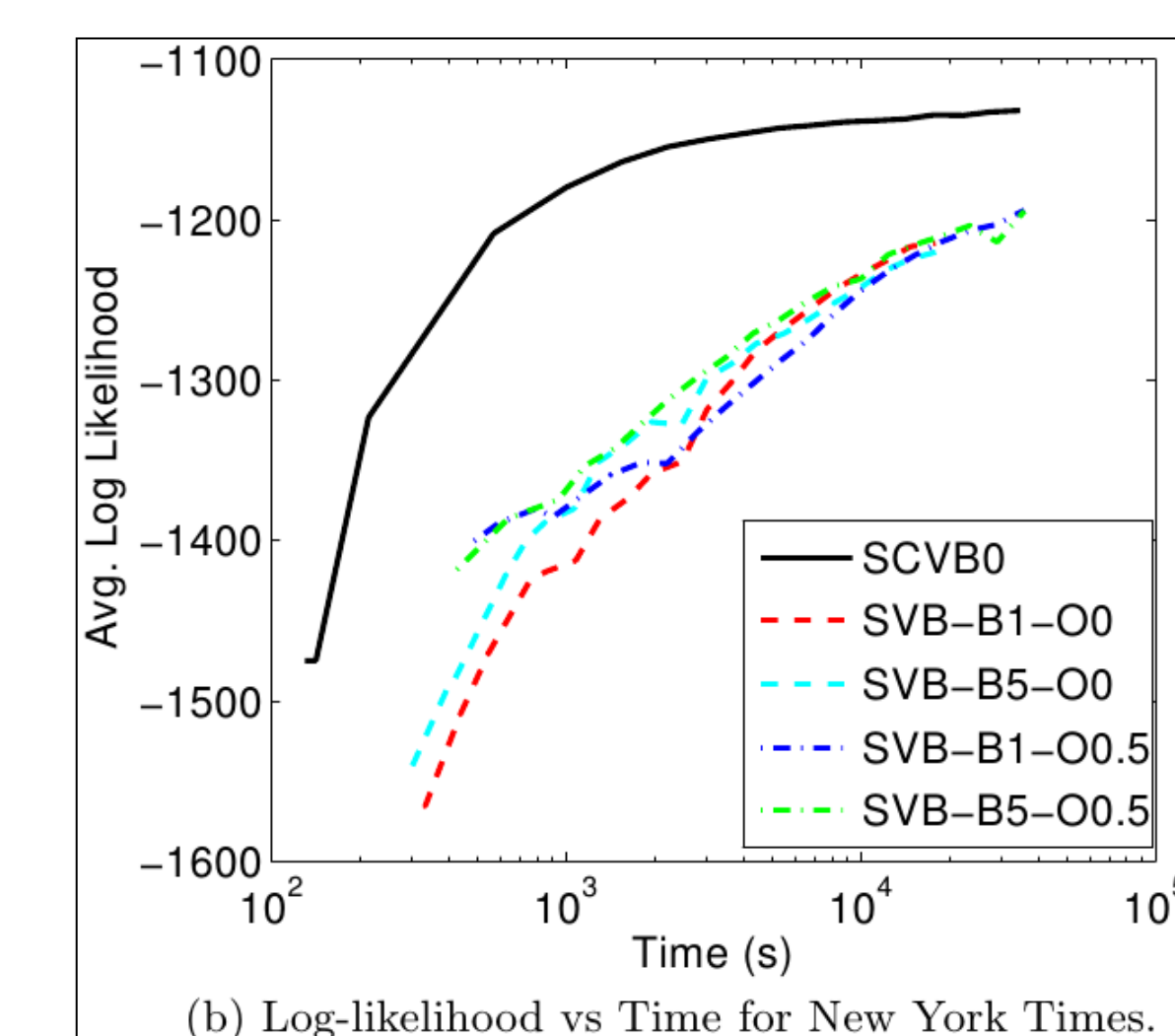
### Large-Scale Experiments



(c) Log-likelihood vs Time for Wikipedia.



(a) Log-likelihood vs Time for PubMed Central.

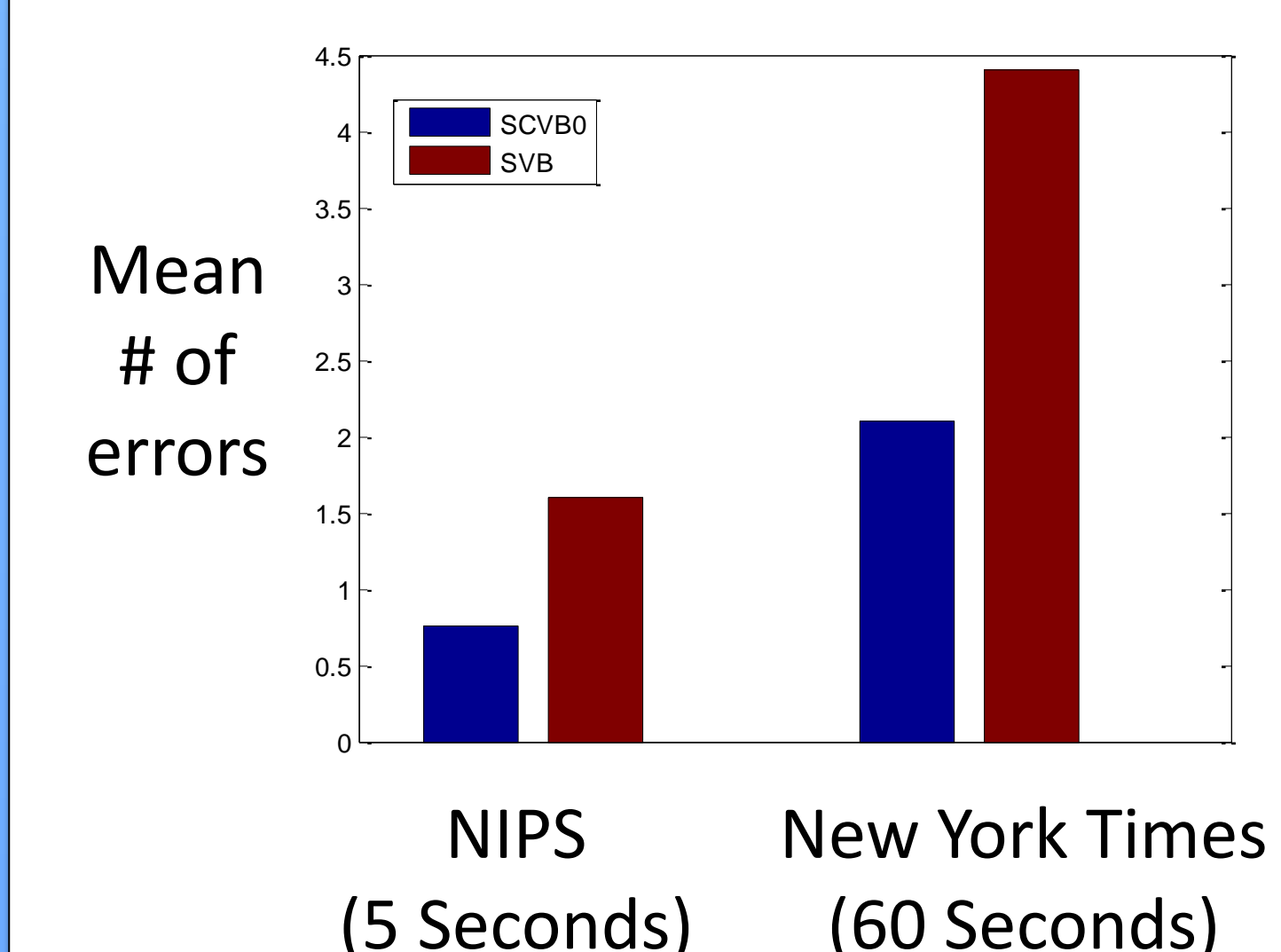


(b) Log-likelihood vs Time for New York Times.

### Small-Scale Experiments

- **Real-time** or near real-time results are important for EDA applications
- Topics on-demand

Here are 20 collections of related words. Some words may not seem to “belong” with the other words. Count the total number of words in each collection that don’t “belong.”



• We introduced stochastic CVB0 for LDA, which is **fast**, **easy** to implement, and **accurate**

• Experimental results show SCVB0 is useful for both **large-scale** and **small-scale** analysis

• Future work: Exploit **sparsity**, **parallelization**, **non-parametric** extensions