



ICML 2017
THE 34TH INTERNATIONAL CONFERENCE ON MACHINE LEARNING



Deep Latent Dirichlet Allocation with Topic-Layer-Adaptive Stochastic Gradient Riemannian MCMC

ICML @ Sydney, August 9, 2017

Yulai Cong*, Bo Chen*, Hongwei Liu*, Mingyuan Zhou#

* National Laboratory of Radar Signal Processing,

* Collaborative Innovation Center of Information
Sensing and Understanding,

Xidian University, Xi'an, Shaanxi, China



西安电子科技大学
XIDIAN UNIVERSITY



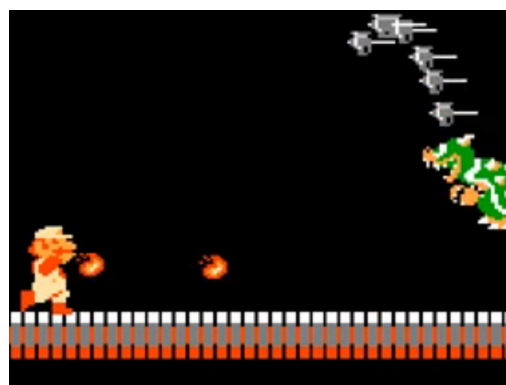
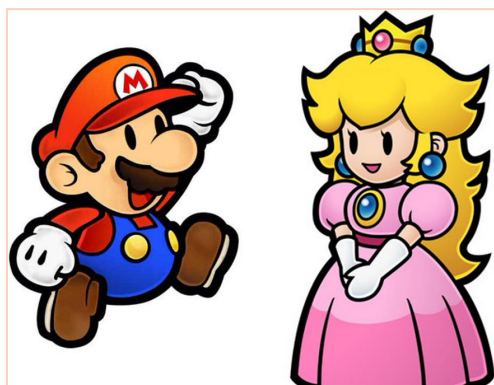
IROM Department, The University of
Texas at Austin, Austin, TX, USA

THE UNIVERSITY OF TEXAS AT AUSTIN





Outline

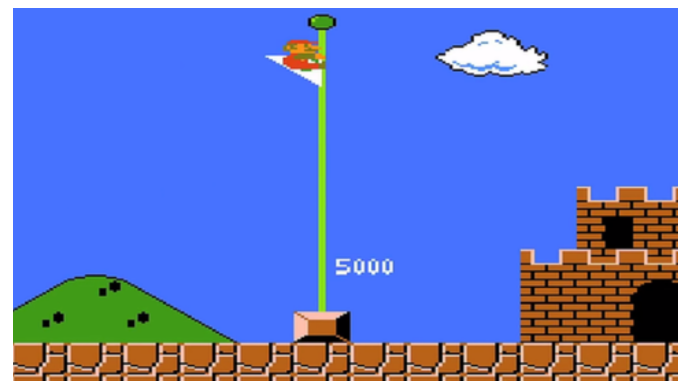


1. Motivation

3. Our Contribution

2. Barriers

4. Experiments





Motivation



Big data have

- ▶ Abundant information
- ▶ Huge volume
- ▶ ...

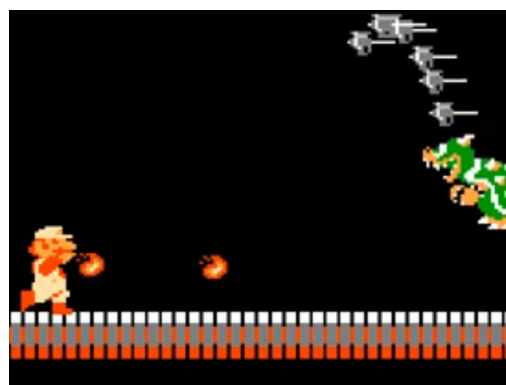
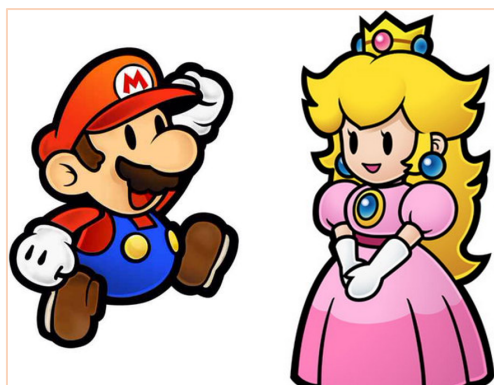
Thus, we prefer

- ▶ Large-capacity models
 - ⌋ Deep latent variable models (LVMs)
 - ⌋ ...
- ▶ Scalable inference methods
 - ⌋ Stochastic Gradient (SG) MCMC
 - ⌋ ...





Outline

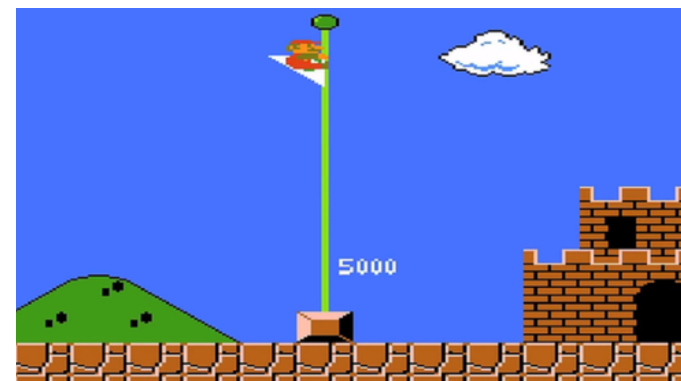


1. Motivation

3. Our Contribution

2. Barriers

4. Experiments

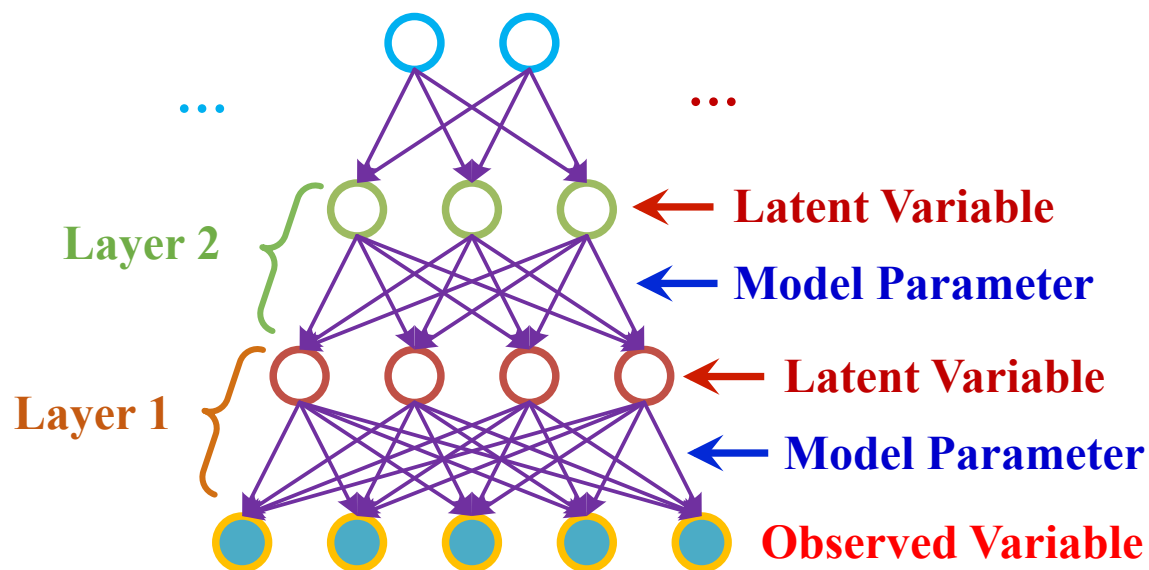




Barriers



To make a deep LVM scalable is challenging.



- ▶ Gradients of model parameters are difficult to compute
- ▶ Different layers, with different statistics, may require different learning rates



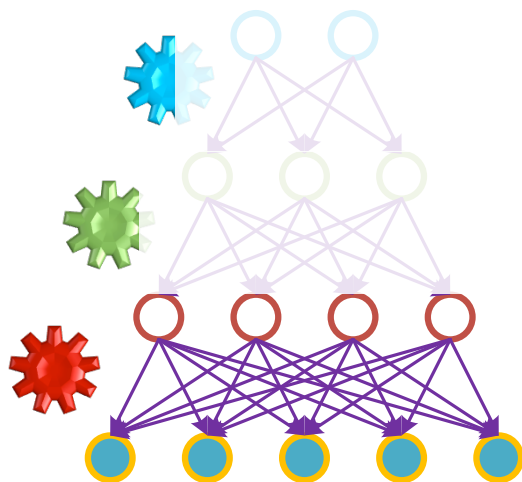


Barriers



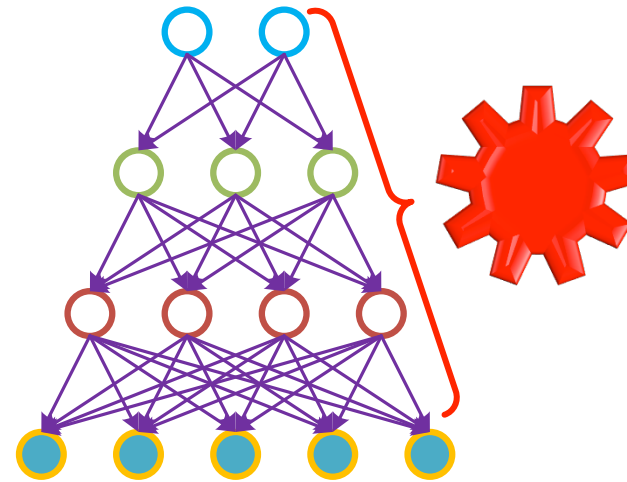
Most existing methods

- ▶ Scalable **but**
- ▶ Greedy layer-wise training
 - ⌈ No communication between layers
- ▶ Shared SGD parameters



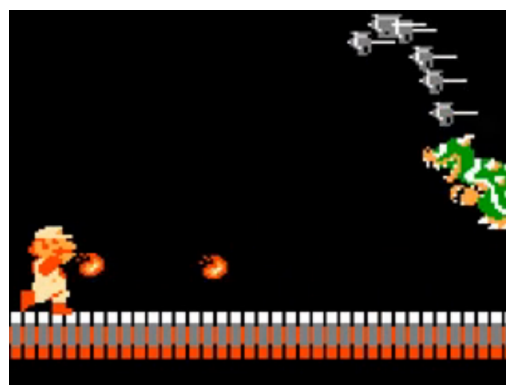
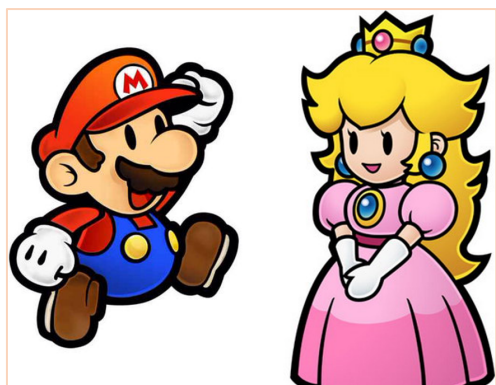
What we want

- ▶ Scalable **and**
- ▶ Principled joint learning
 - ⌈ Communication & feedback between layers
- ▶ Adaptive SGD parameters





Outline

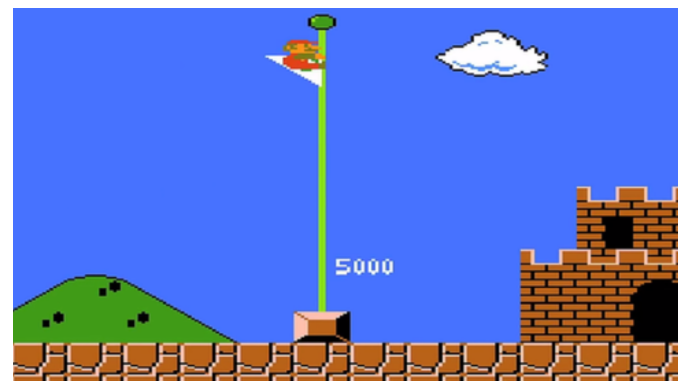


1. Motivation

3. Our Contribution

2. Barriers

4. Experiments





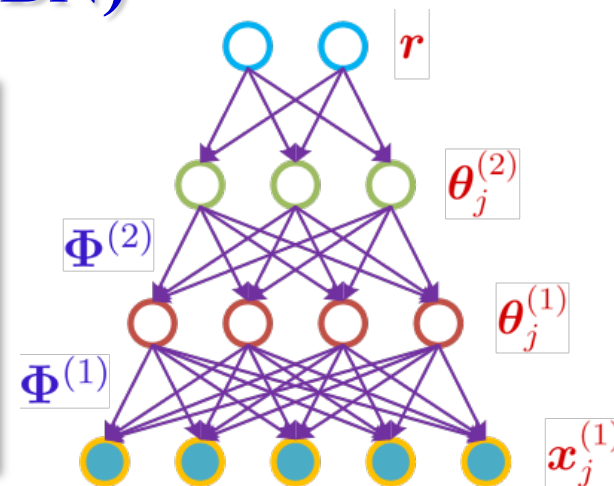
Our Contribution



We develop a **principled joint SG-MCMC** for a special deep LVM, namely the **Poisson gamma belief network (PGBN)**.

□ Poisson Gamma Belief Network (PGBN)

$$\begin{aligned} \theta_j^{(L)} &\sim \text{Gam} \left(\mathbf{r}, 1/c_j^{(L+1)} \right), \\ &\dots \\ \theta_j^{(l)} &\sim \text{Gam} \left(\Phi^{(l+1)} \theta_j^{(l+1)}, 1/c_j^{(l+1)} \right), \\ &\dots \\ \mathbf{x}_j^{(1)} &\sim \text{Pois} \left(\Phi^{(1)} \theta_j^{(1)} \right), \quad \theta_j^{(1)} \sim \text{Gam} \left(\Phi^{(2)} \theta_j^{(2)}, \frac{p_j^{(2)}}{1-p_j^{(2)}} \right), \end{aligned}$$



Priors:

$$\phi_k^{(l)} \sim \text{Dir}(\eta^{(l)} \mathbf{1}_{K_{l-1}})$$

$$\mathbf{r} \sim \text{Gam}(\gamma_0/K_L, 1/c_0)$$

$$p_j^{(2)} \sim \text{Beta}(a_0, b_0)$$

$$c_j^{(l)} \sim \text{Gam}(e_0, 1/f_0)$$

Gibbs

Not scalable





Preliminary for SG-MCMC



For the interested **batch posterior** $p(\mathbf{z} | X) \propto e^{-H(\mathbf{z})}$, one may get **posterior samples** using the **mini-batch** update rule as

$$\mathbf{z}_{t+1} = \mathbf{z}_t + \varepsilon_t \left\{ - [\mathbf{D}(\mathbf{z}_t) + \mathbf{Q}(\mathbf{z}_t)] \nabla \tilde{H}(\mathbf{z}_t) + \Gamma(\mathbf{z}_t) \right\} + \mathcal{N}\left(\mathbf{0}, \varepsilon_t [2\mathbf{D}(\mathbf{z}_t) - \varepsilon_t \hat{\mathbf{B}}_t]\right),$$

$$\nabla \tilde{H}(\mathbf{z}) = \nabla \left[-\ln p(\mathbf{z}) - \rho \sum_{\mathbf{x} \in \tilde{X}} \ln p(\mathbf{x} | \mathbf{z}) \right]$$

$\mathbf{D}(\mathbf{z})$: positive semi-definite;

$\mathbf{Q}(\mathbf{z})$: skew-symmetric;

$$\Gamma_i(\mathbf{z}) = \sum_j \frac{\partial}{\partial z_j} [\mathbf{D}_{ij}(\mathbf{z}) + \mathbf{Q}_{ij}(\mathbf{z})]$$

$\hat{\mathbf{B}}_t$: SG noise variance estimate.





SG-MCMC for PGBN



SG-MCMC & PGBN

▶ $\mathbf{z} = \{ \Phi^{(1)}, \dots, \Phi^{(L)}, \mathbf{r} \}$

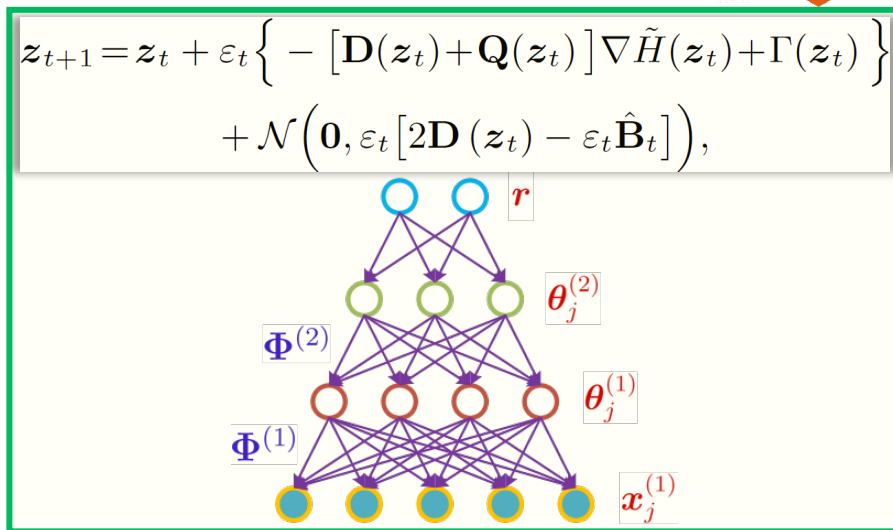
For simplicity, we use

▶ $\mathbf{Q}(\mathbf{z}) = \mathbf{0}, \hat{\mathbf{B}}_t = \mathbf{0}$

▶ $\mathbf{D}(\mathbf{z}) = \mathbf{G}(\mathbf{z})^{-1}$

⤴ $\mathbf{G}(\mathbf{z})$: Fisher information matrix (FIM)

$$\mathbf{z}_{t+1} = \mathbf{z}_t + \varepsilon_t \left\{ -\mathbf{G}(\mathbf{z}_t)^{-1} \nabla \tilde{H}(\mathbf{z}_t) + \Gamma(\mathbf{z}_t) \right\} + \mathcal{N}(\mathbf{0}, 2\varepsilon_t \mathbf{G}(\mathbf{z}_t)^{-1})$$



We need

▶ $\mathbf{G}(\mathbf{z})^{-1}$

▶ $\nabla \tilde{H}(\mathbf{z})$

However

▶ $\mathbf{G}(\mathbf{z})$ of PGBN is intractable

▶ $\mathbf{G}(\mathbf{z})^{-1}$ may be expensive





Calculate FIM



What we find: An alternative representation of the PGBN makes it straightforward to calculate $\mathbf{G}(z)$.

Deep Latent Dirichlet Allocation (DLDA)

Techniques:

► Data augmentation

► Marginalization

□ $q_j^{(1)} := 1$

□ $q_j^{(l+1)} = \ln \left(1 + q_j^{(l)} / c_j^{(l+1)} \right)$

□ $p_j^{(l)} := 1 - e^{-q_j^{(l)}}$

□ $\tilde{p} := q_{\cdot}^{(L+1)} / (c_0 + q_{\cdot}^{(L+1)})$

$$\begin{aligned}
x_{k\cdot}^{(L+1)} &\sim \text{Log}(\tilde{p}), K_L \sim \text{Pois}[-\gamma_0 \ln(1 - \tilde{p})], \\
X^{(L+1)} &= \sum_{k=1}^{K_L} x_{k\cdot}^{(L+1)} \delta_{\phi_k^{(L)}}, \\
(x_{vj}^{(L+1)})_j &\sim \text{Mult} \left[x_{v\cdot}^{(L+1)}, (q_j^{(L+1)})_j / q_{\cdot}^{(L+1)} \right], \\
m_{vj}^{(L)(L+1)} &\sim \text{SumLog}(x_{vj}^{(L+1)}, p_j^{(L+1)}), \\
&\dots \\
x_{vj}^{(l)} &= \sum_{k=1}^{K_l} x_{v kj}^{(l)}, (x_{v kj}^{(l)})_v \sim \text{Mult} \left(m_{kj}^{(l)(l+1)}, \phi_k^{(l)} \right), \\
m_{vj}^{(l-1)(l)} &\sim \text{SumLog}(x_{vj}^{(l)}, p_j^{(l)}), \\
&\dots \\
x_{vj}^{(1)} &= \sum_{k=1}^{K_1} x_{v kj}^{(1)}, (x_{v kj}^{(1)})_v \sim \text{Mult} \left(m_{kj}^{(1)(2)}, \phi_k^{(1)} \right).
\end{aligned}$$





Calculate FIM



The alternative DLDA representation \rightarrow
A *block-diagonal* and thus *easily inverted* FIM as

$$\mathbf{G}(\mathbf{z}) = \text{diag} \left[\mathbf{I} \left(\phi_1^{(1)} \right), \dots, \mathbf{I} \left(\phi_{K_L}^{(L)} \right), \mathbf{I}(\mathbf{r}) \right]$$

Properties

- ▶ Analytical & Practical
- ▶ Similar to the Hessian matrix in optimization
- ▶ Principled joint inference of the PGBN (DLDA)
- ▶ A separate puzzle for each $\phi_k^{(l)}$ and \mathbf{r}

$$\mathbf{z}_{t+1} = \mathbf{z}_t + \varepsilon_t \left\{ -\mathbf{G}(\mathbf{z}_t)^{-1} \nabla \tilde{H}(\mathbf{z}_t) + \Gamma(\mathbf{z}_t) \right\} + \mathcal{N}(\mathbf{0}, 2\varepsilon_t \mathbf{G}(\mathbf{z}_t)^{-1})$$





Efficient Inference for Topics $\phi_k^{(l)}$



Reduced-Mean Inference on the Probability Simplex

For the batch posterior $(\phi_k | -) \sim \text{Dir}(x_{1k\cdot} + \eta, \dots, x_{V_k\cdot} + \eta)$ on the probability simplex, we derive a *mini-batch* update rule as

$$(\phi_k)_{t+1} = \left[(\phi_k)_t + \frac{\varepsilon_t}{M_k} \left[(\rho \tilde{x}_{:k\cdot} + \eta) - (\rho \tilde{x}_{\cdot k\cdot} + \eta V) (\phi_k)_t \right] + \mathcal{N}(\mathbf{0}, \frac{2\varepsilon_t}{M_k} \text{diag}(\phi_k)_t) \right]_{\Delta}$$

where $\tilde{x}_{vk\cdot}$ comes from mini-batches.

Techniques

► Reduced-mean parameterization φ_k of ϕ_k , namely

$$\phi_k = \left((\varphi_k)^T, 1 - \sum_v \varphi_{vk} \right)^T$$

► Fast simulation method in [1]

[1] Cong, Y., Chen, B., and Zhou, M. Fast simulation of hyperplane-truncated multivariate normal distributions. Bayesian Analysis Advance Publication, 2017.





The Whole Picture



□ Topic-Layer-Adaptive Stochastic Gradient Riemannian (TLASGR) MCMC

Algorithm 1 TLASGR MCMC for DLDA (PGBN).

Input: Data mini-batches;

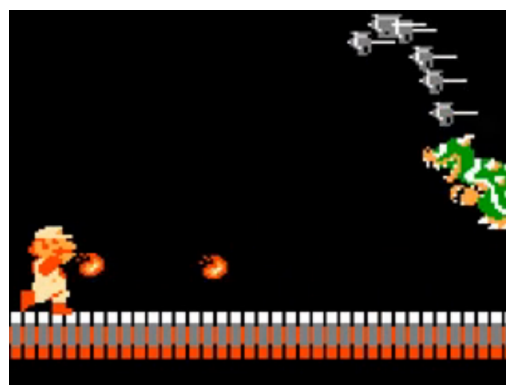
Output: Global parameters of DLDA (PGBN).

```
1: for  $t = 1, 2, \dots$  do
2:   /* Collect local information
3:   Upward-downward Gibbs sampling (Zhou et al., 2016a) on
   the  $t^{\text{th}}$  mini-batch for  $\tilde{\mathbf{x}}_{:k.}, \tilde{\mathbf{x}}_{\cdot k.}, \tilde{\mathbf{x}}_{: \cdot}^{(L+1)}$ , and  $\tilde{\mathbf{q}}_{\cdot}^{(L+1)}$ ;
4:   /* Update global parameters
5:   for  $l = 1, \dots, L$  and  $k = 1, \dots, K_l$  do
6:     Update  $M_k^{(l)}$  with (18); then  $\phi_k^{(l)}$  with (15);
7:   end for
8:   Update  $M^{(L+1)}$  with (19) and then  $\mathbf{r}$  with (17).
9: end for
```





Outline

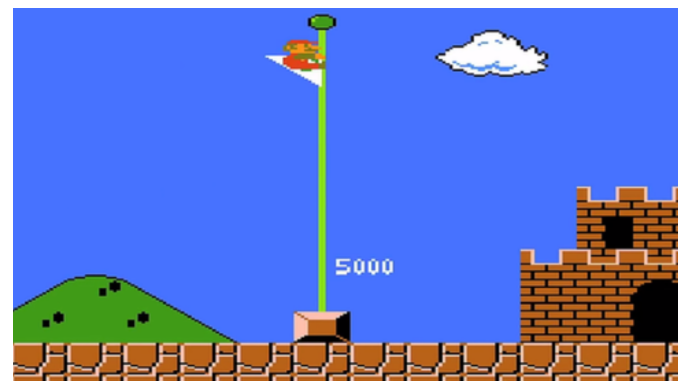


1. Motivation

3. Our Contribution

2. Barriers

4. Experiments





Experiments



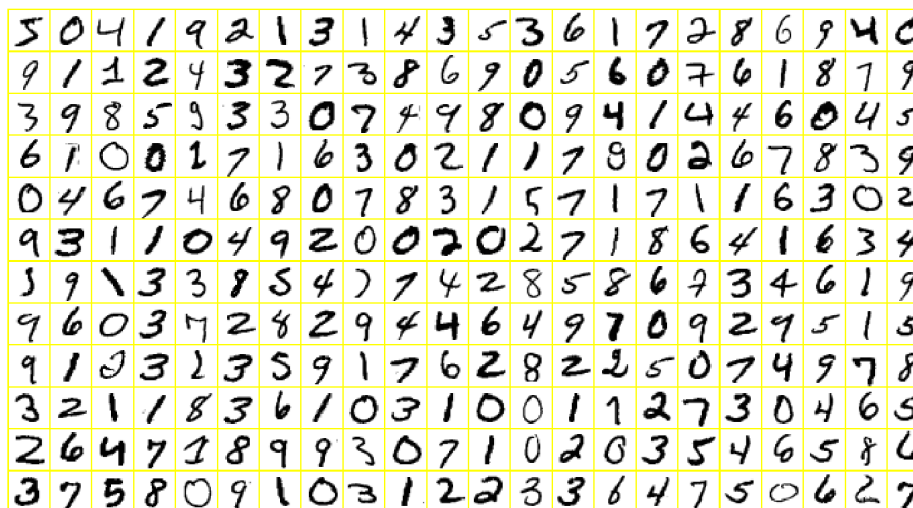
Three benchmark corpora

- ▶ 20Newsgroups (20News)
- ▶ Reuters Corpus Volume I (RCV1)
- ▶ Wikipedia (Wiki)

	20News	RCV1	Wiki
Vocabulary	2,000	10,000	7,702
Training	11,315	794,414	≈10M
Test	7,531	10,000	1,000

MNIST digits

- ▶ For illustration
- ▶ Size: 28×28
- ▶ Training: 60000
- ▶ Test: 10000





Performance: Perplexity



▶ **Smaller is better**

▶ **Bold: Top two**

DLDA & Gibbs

▶ **Best results**

▶ **Not scalable**

DLDA & TLASGR

▶ **Second best**

‡ **Comparable to Gibbs**

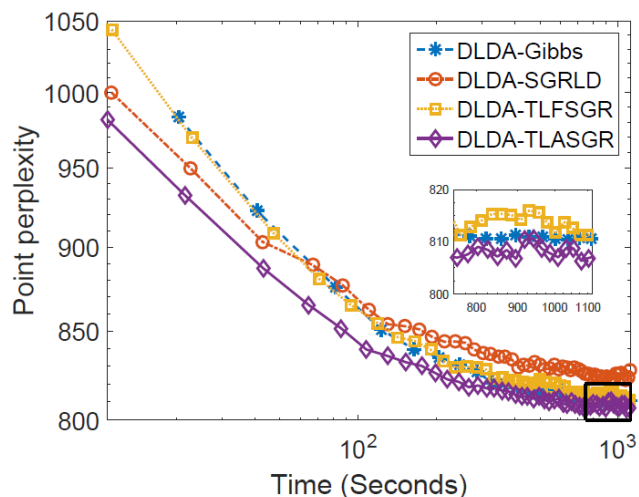
▶ **Scalable**

Model	Method	Size	20 News	RCV1	Wiki
DLDA	TLASGR	128-64-32	757	815	786
DLDA	TLASGR	128-64	758	817	787
DLDA	TLASGR	128	770	823	802
DLDA	TLFSGR	128-64-32	760	817	789
DLDA	TLFSGR	128-64	759	819	791
DLDA	TLFSGR	128	772	829	804
DLDA	SGRLD	128-64-32	775	827	792
DLDA	SGRLD	128-64	770	823	792
DLDA	SGRLD	128	777	829	803
DLDA	Gibbs	128-64-32	752	802	—
DLDA	Gibbs	128-64	754	804	—
DLDA	Gibbs	128	768	818	—
DPFM	SVI	128-64	818	961	791
DPFM	MCMC	128-64	780	908	783
DPFA-SBN	SGNHT	128-64-32	827	1143	876
DPFA-RBM	SGNHT	128-64-32	896	920	942
nHDP	SVI	(10,10,5)	889	1041	932
LDA	Gibbs	128	893	1179	1059
FTM	Gibbs	128	887	1155	991
RSM	CD5	128	877	1171	1001

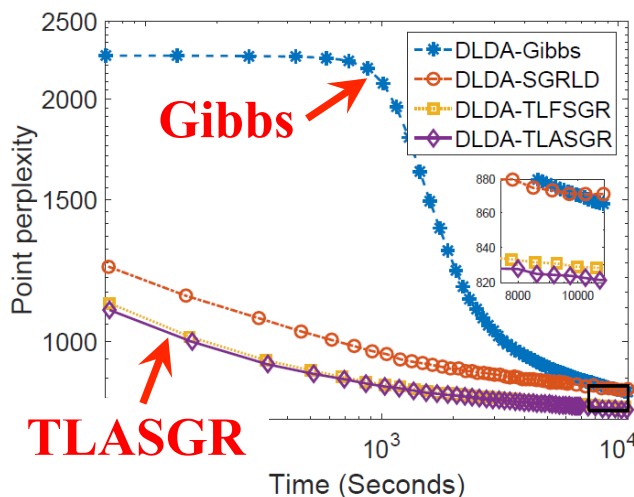




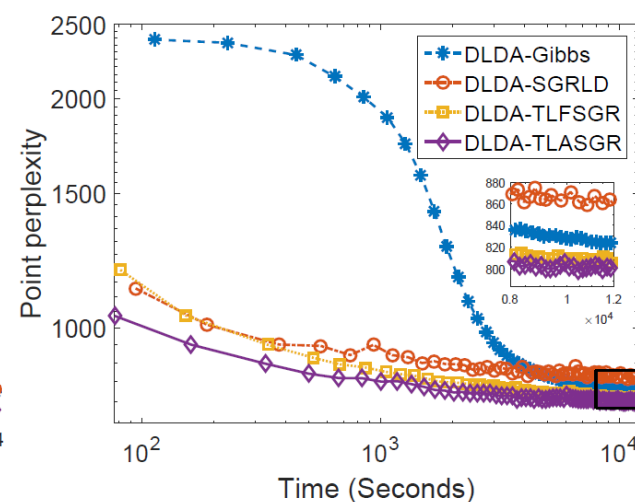
Performance: Scalability



(a) A single-layer DLDA on 20News



(b) DLDA of size 128-64 on RCV1



(c) DLDA of size 128-64 on Wiki

20News
Train: 11K

RCV1
Train: 0.8M

Wiki
Train: 10M

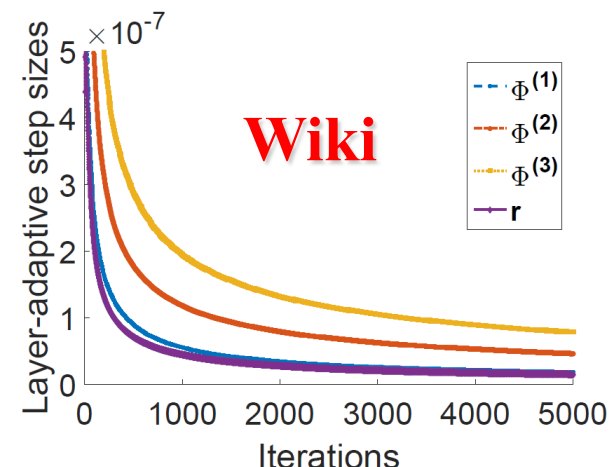
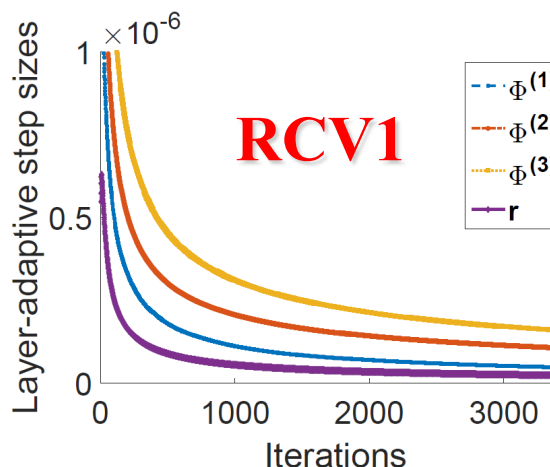
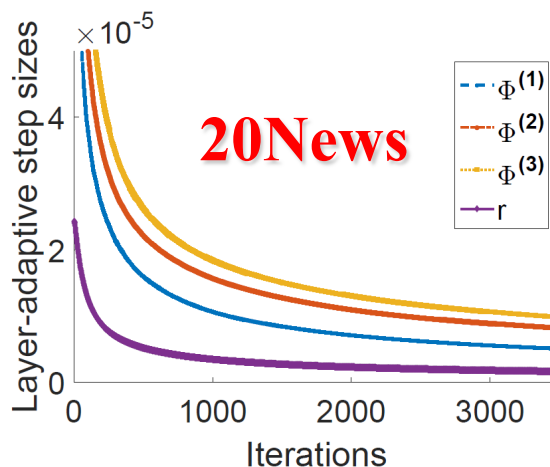




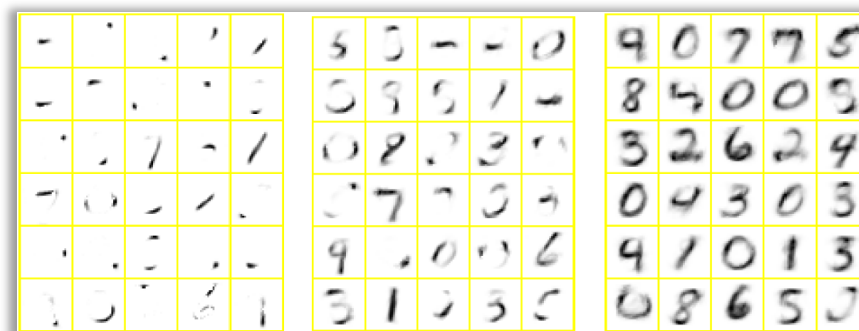
Benefits of Joint Learning



► Topic-layer-adaptive learning rates

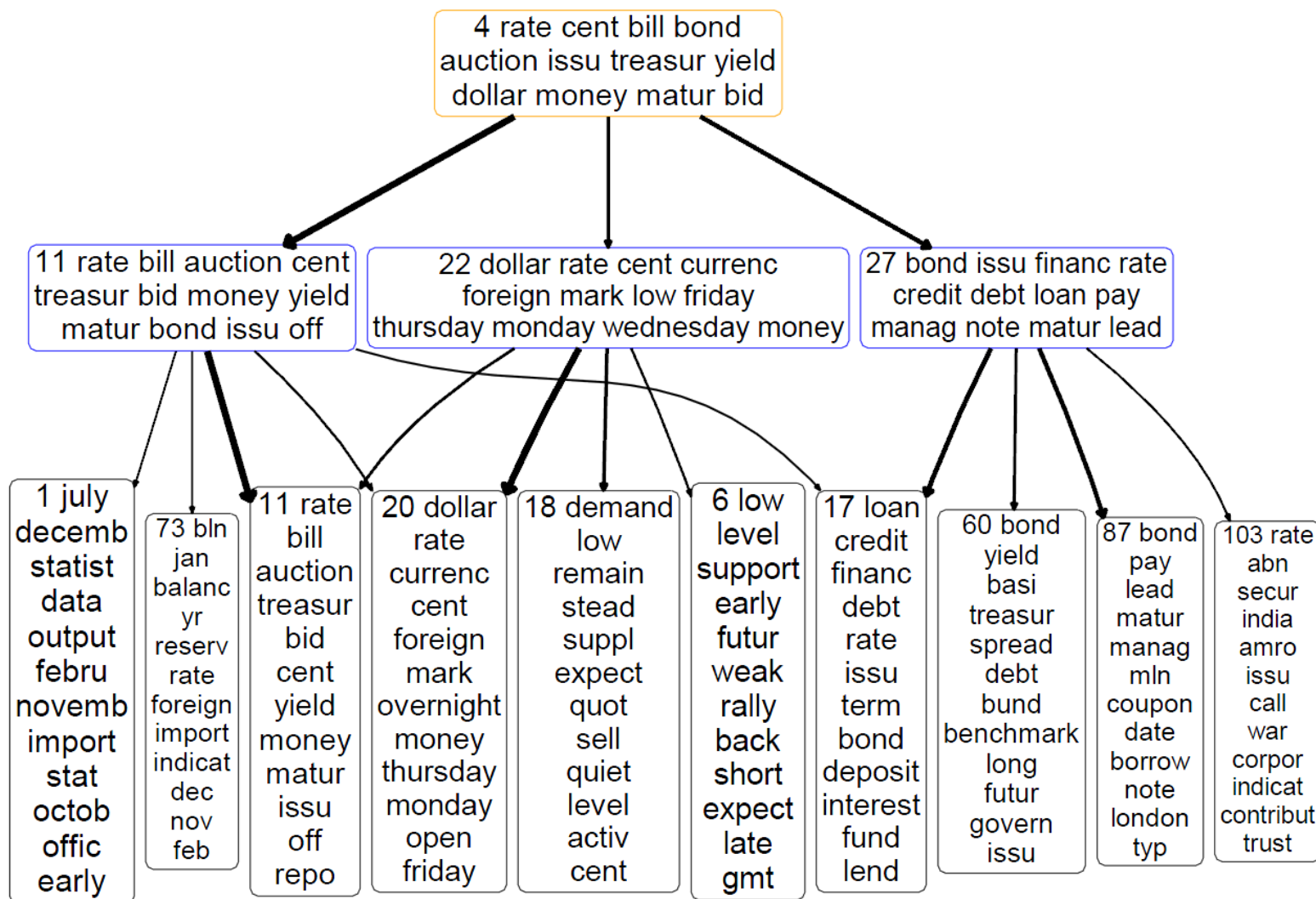


► Better information propagation



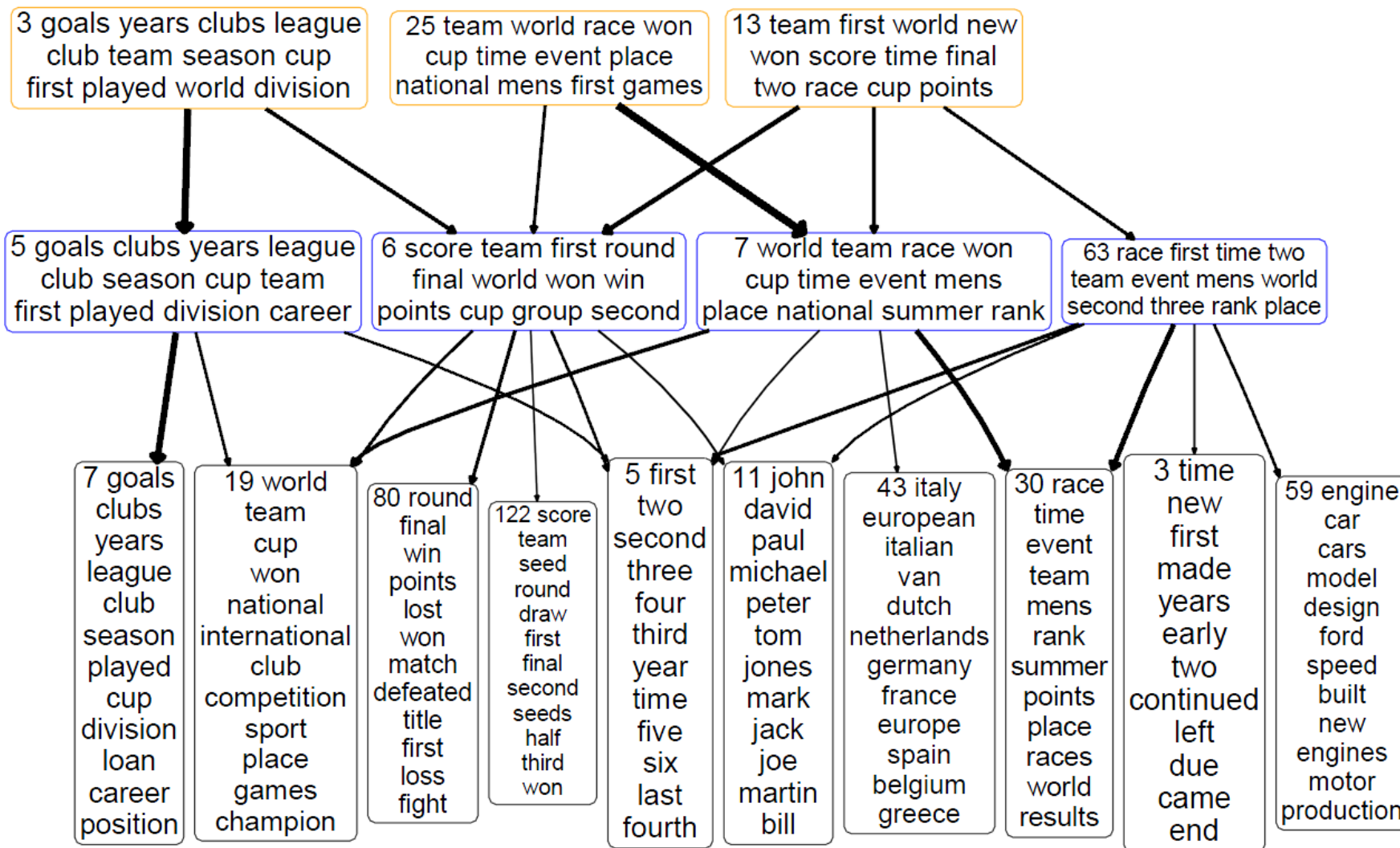


Example Topics: RCV1





Example Topics: Wiki

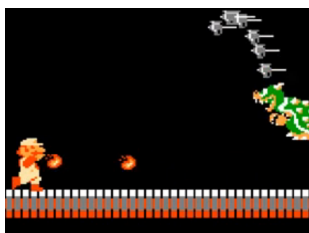




Conclusions



Scalable inference of deep LVMs for big data

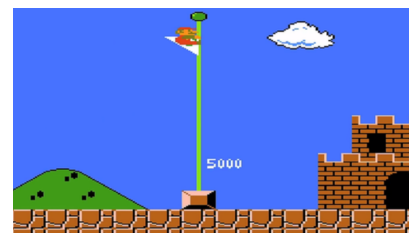


For the PGBN, we develop the principled joint **TLASGR** via newly introduced

- ▶ An alternative representation DLDA
- ▶ Analytical and block-diagonal FIM
- ▶ Reduced-mean parameterization on probability simplex



- ▶ Gradients are difficult to calculate
- ▶ Joint learning of all layers of deep models



Experimentally, **TLASGR**

- ▶ Performance comparable to Gibbs
- ▶ Scalable to huge data
- ▶ Principle joint learning
 - ⌞ Topic layer adaptive learning rates
 - ⌞ Better information propagation
- ▶ Interpretable topics





THANK YOU!

