

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

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

* National Laboratory of Radar Signal Processing, Xidian University, Xi'an, Shaanxi, China
IROM Department, The University of Texas at Austin, Austin, TX, USA



Motivation

Big data prefer
 ➤ large-capacity models like deep latent variable models (LVMs)
 ➤ scalable inference methods like SG-MCMC

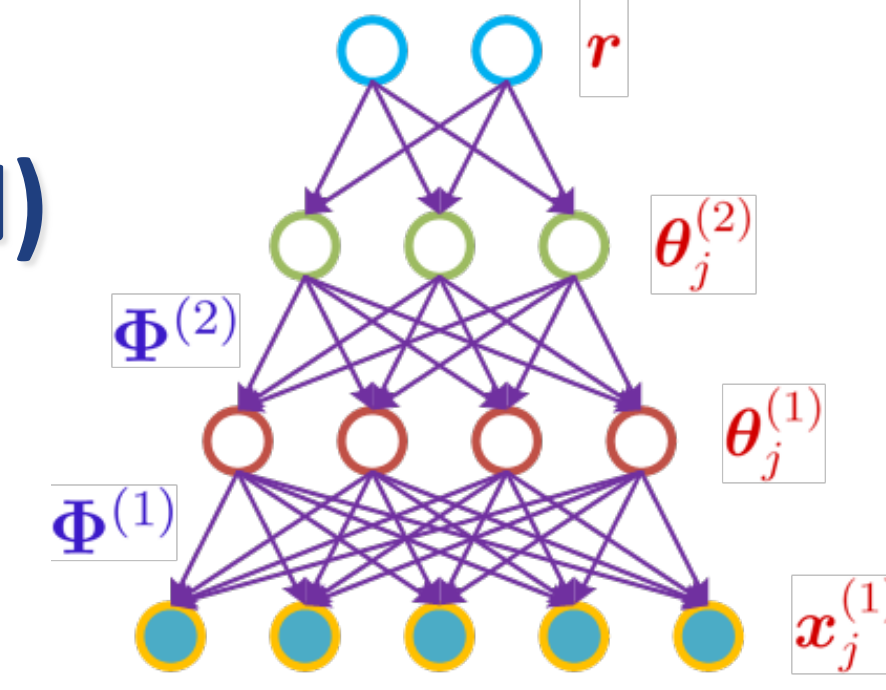
However, **to make a deep LVM scalable is challenging**, because
 ➤ gradients of model parameters are difficult to compute
 ➤ **different layers may be suitable for different learning rates**

Most **existing methods** adopt greedy layer-wise training, which
 ➤ usually uses a shared learning rates across all layers
 ➤ **has no communication between different layers**

We develop a scalable SG-MCMC algorithm for deep latent Dirichlet allocation that jointly learns all hidden layers.

Background

Poisson Gamma Belief Network (PGBN)



$$\theta_j^{(L)} \sim \text{Gam}(r, 1/c_j^{(L+1)}),$$

$$\theta_j^{(l)} \sim \text{Gam}(\Phi^{(l+1)}\theta_j^{(l+1)}, 1/c_j^{(l+1)}),$$

$$x_j^{(1)} \sim \text{Pois}(\Phi^{(1)}\theta_j^{(1)}), \theta_j^{(1)} \sim \text{Gam}(\Phi^{(2)}\theta_j^{(2)}, \frac{p_j^{(2)}}{1-p_j^{(2)}}),$$

Stochastic Gradient MCMC

For unknown z obeying $p(z|X) \propto e^{-H(z)}$, one has a mini-batch update rule as

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

where $\mathbf{D}(z)$ is a positive semidefinite matrix, $\mathbf{Q}(z)$ is a skew-symmetric matrix, $\Gamma_i(z) = \sum_j \frac{\partial}{\partial z_j} [\mathbf{D}_{ij}(z) + \mathbf{Q}_{ij}(z)]$, $\nabla \tilde{H}(z) = \nabla [-\ln p(z) - \rho \sum_{x \in \tilde{X}} \ln p(x|z)]$ is estimated with mini-batches, and $\hat{\mathbf{B}}_t$ is the SG noise variance estimate.

Contributions

To develop principled SG-MCMC for the PGBN, with $z = \{\Phi^{(1)}, \dots, \Phi^{(L)}, r\}$ and simple settings $\mathbf{D}(z) = \mathbf{G}(z)^{-1}$, $\mathbf{Q}(z) = \mathbf{0}$, and $\hat{\mathbf{B}}_t = \mathbf{0}$, all we need are (a) the Fisher information matrix (FIM) $\mathbf{G}(z)$ and (b) $\nabla \tilde{H}(z)$.

Directly computing the FIM of PGBN is intractable. But an alternative view of the PGBN makes it straightforward.

Deep Latent Dirichlet Allocation (DLDA)

Exploiting data augmentation and marginalization techniques, one may re-express the hierarchical model of the PGBN as DLDA as

$$x_k^{(L+1)} \sim \text{Log}(\tilde{p}), K_L \sim \text{Pois}[-\gamma_0 \ln(1 - \tilde{p})], X^{(L+1)} = \sum_{k=1}^{K_L} x_k^{(L+1)} \delta_{\phi_k^{(L)}},$$

$$(x_{vj}^{(L+1)})_j \sim \text{Mult} \left[x_{v \cdot}^{(L+1)}, (q_j^{(L+1)})_j / q^{(L+1)} \right], m_{vj}^{(L+1)} \sim \text{SumLog}(x_{vj}^{(L+1)}, p_j^{(L+1)}),$$

...

$$x_{vj}^{(l)} = \sum_{k=1}^{K_l} x_{v_k j}^{(l)}, (x_{v_k j}^{(l)})_v \sim \text{Mult}(m_{kj}^{(l+1)}, \phi_k^{(l)}), m_{vj}^{(l-1)} \sim \text{SumLog}(x_{vj}^{(l)}, p_j^{(l)}),$$

...

$$x_{vj}^{(1)} = \sum_{k=1}^{K_1} x_{v_k j}^{(1)}, (x_{v_k j}^{(1)})_v \sim \text{Mult}(m_{kj}^{(1)}, \phi_k^{(1)}).$$

where $q_j^{(1)} := 1$, $q_j^{(l+1)} = \ln(1 + q_j^{(l)}/c_j^{(l+1)})$, $p_j^{(l)} := 1 - e^{-q_j^{(l)}}$, and $\tilde{p} := q^{(L+1)}/(c_0 + q^{(L+1)})$.

Analytical and Practical Fisher Information Matrix

Under the alternative DLDA representation, one may readily derive a **block-diagonal** and thus **easily-inversed** Fisher information matrix (FIM) as

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

• Note the FIM, playing a similar role as the Hessian matrix in optimization, enables principled joint inference of the PGBN (DLDA).

Reduced-Mean Inference on the Probability Simplex

For the interested batch posterior $(\phi_k | -) \sim \text{Dir}(x_{1:k} + \eta, \dots, x_{V:k} + \eta)$ on the probability simplex, one may have the mini-batch-based inference as

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

which is derived with the reduced-mean parameterization φ_k of ϕ_k , namely

$$\phi_k = \left((\varphi_k)^T, 1 - \sum_v \varphi_{vk} \right)^T, \text{ and Theorem 2 of [1].}$$

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

Topic-Layer-Adaptive Stochastic Gradient Riemannian (TLASGR)

Algorithm 1 TLASGR MCMC for DLDA (PGBN).

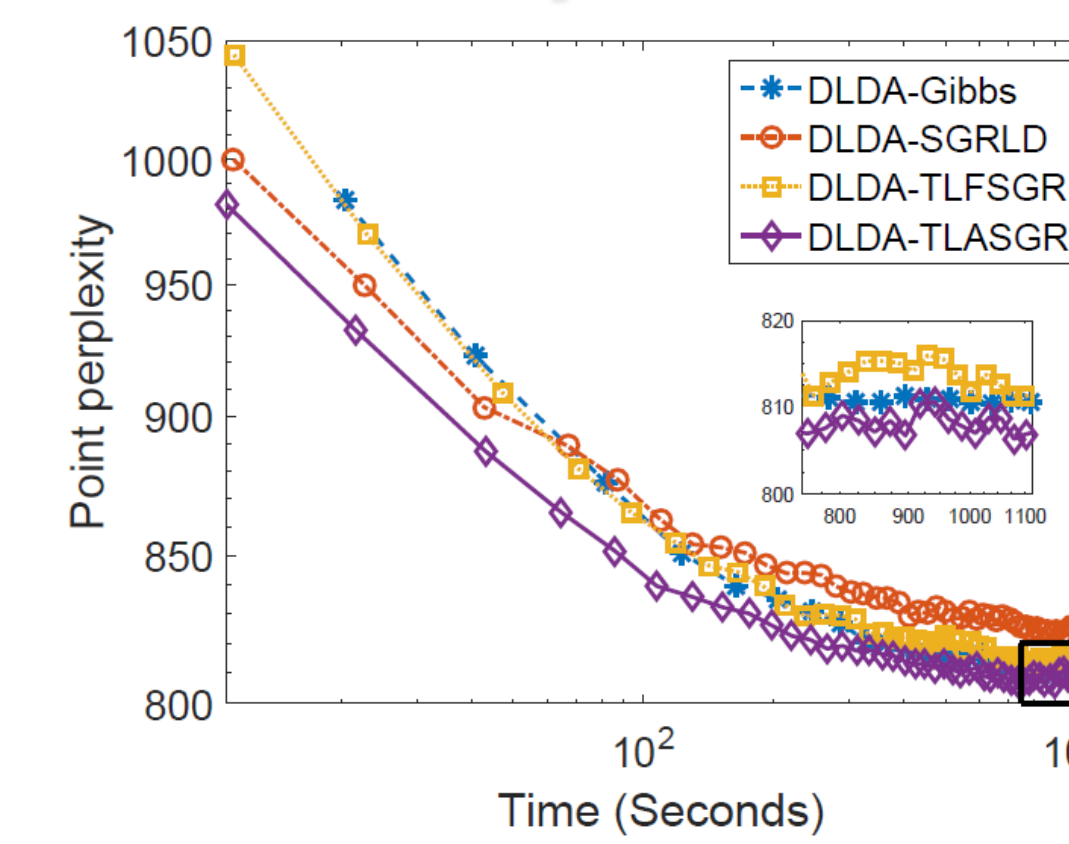
Input: Data mini-batches;
Output: Global parameters of DLDA (PGBN).
 1: **for** $t = 1, 2, \dots$ **do**
 2: /* Collect local information
 3: Upward-downward Gibbs sampling (Zhou et al., 2016a) on the t^{th} mini-batch for $\tilde{x}_{:k}, \tilde{x}_{\cdot k}, \tilde{x}_{:}^{(L+1)}$, and $\tilde{q}^{(L+1)}$;
 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 r with (17).
 9: **end for**

Experiment Results

Perplexity Performance

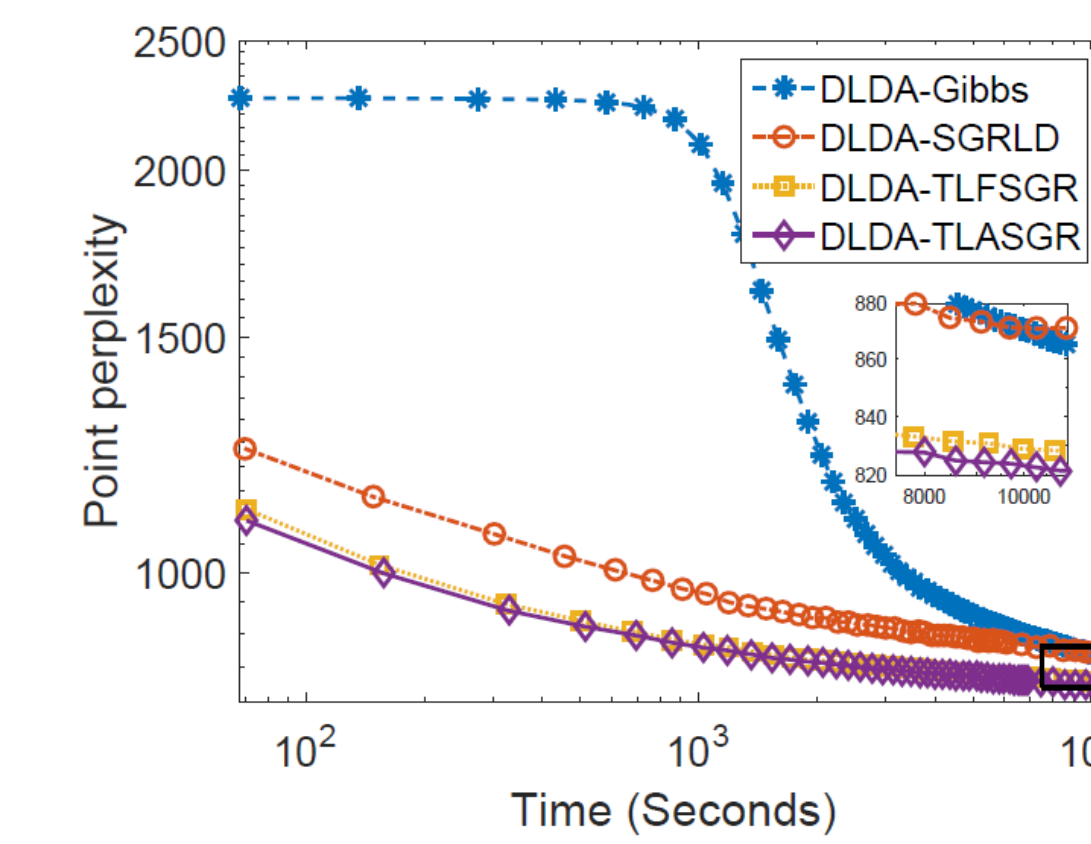
Model	Method	Size	20 News	RCV1	Wiki
DLDA	TLASGR	128-64-32	757	815	786
DLDA	TLASGR	128-64	758	817	787
DLDA	TLASGR	128	770	823	802
DLDA	TLFSGR	128-64-32	760	817	789
DLDA	TLFSGR	128-64	759	819	791
DLDA	TLFSGR	128	772	829	804
DLDA	SGRLD	128-64-32	775	827	792
DLDA	SGRLD	128-64	770	823	792
DLDA	SGRLD	128	777	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

Scalability

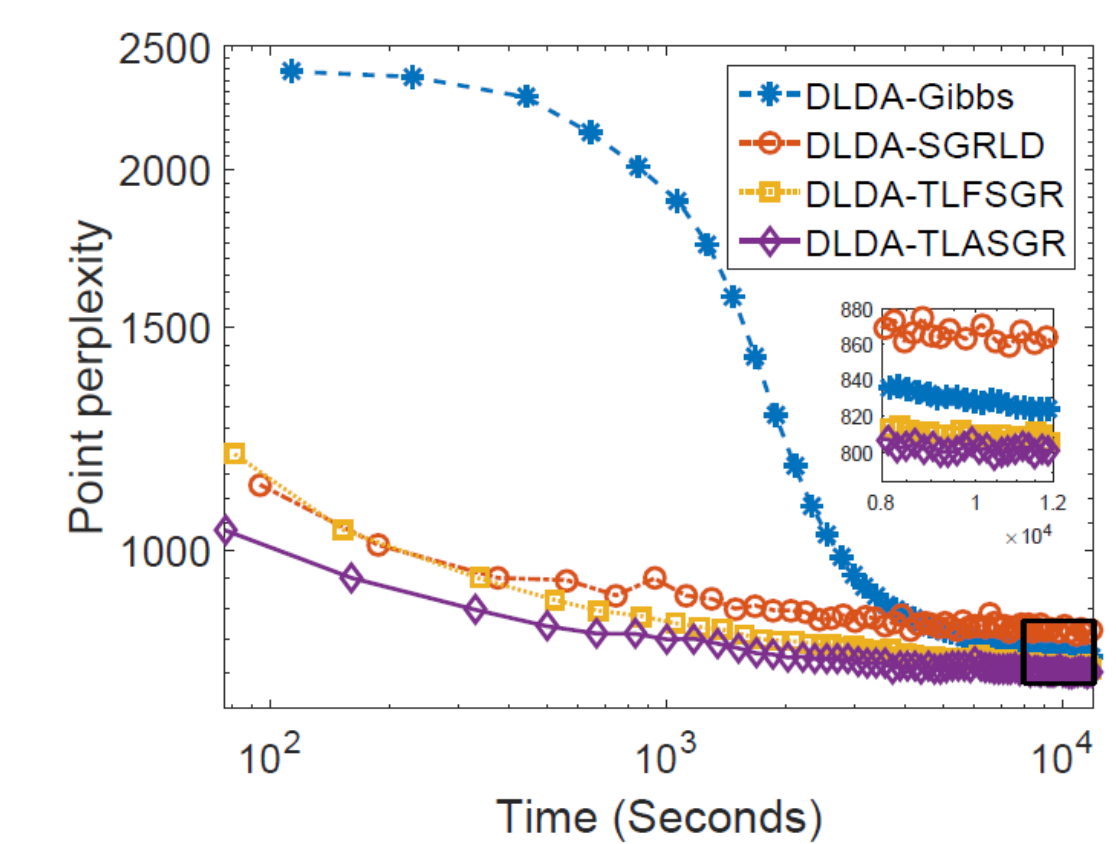


(a) A single-layer DLDA on 20News

Training Docs: 20News 11K, RCV1 0.8M, Wiki 10M

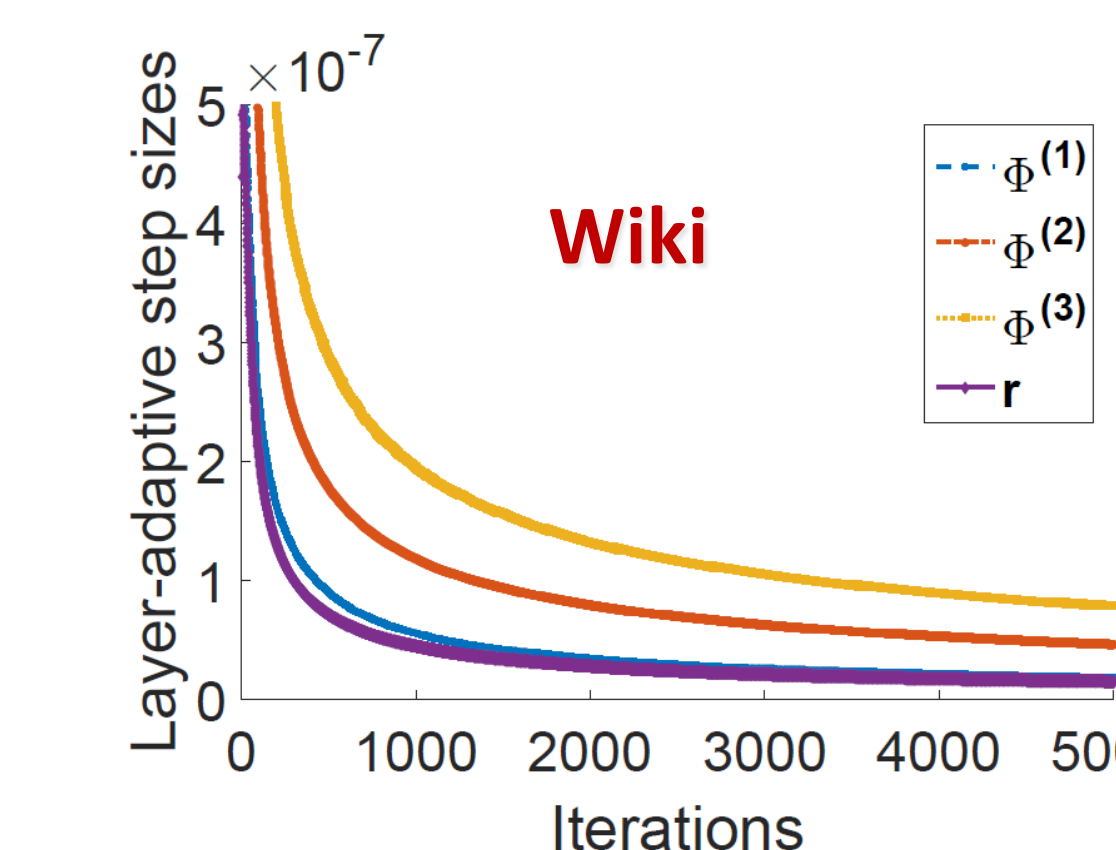
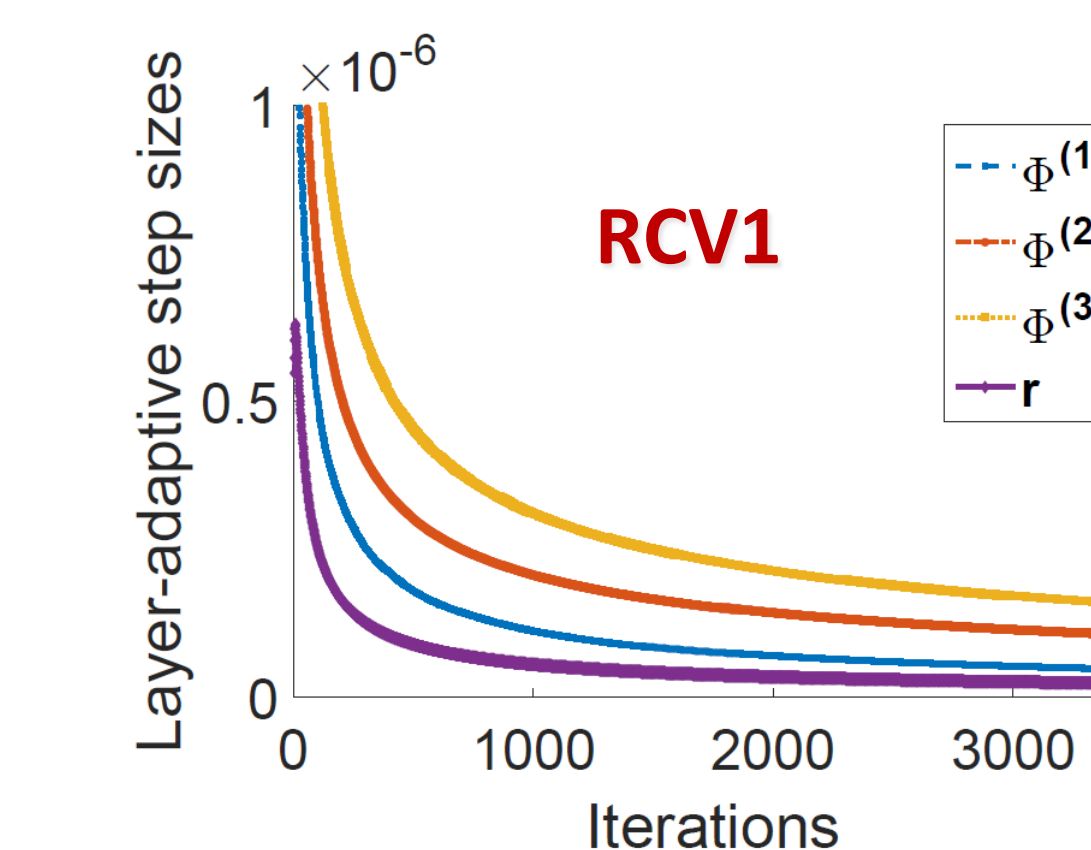
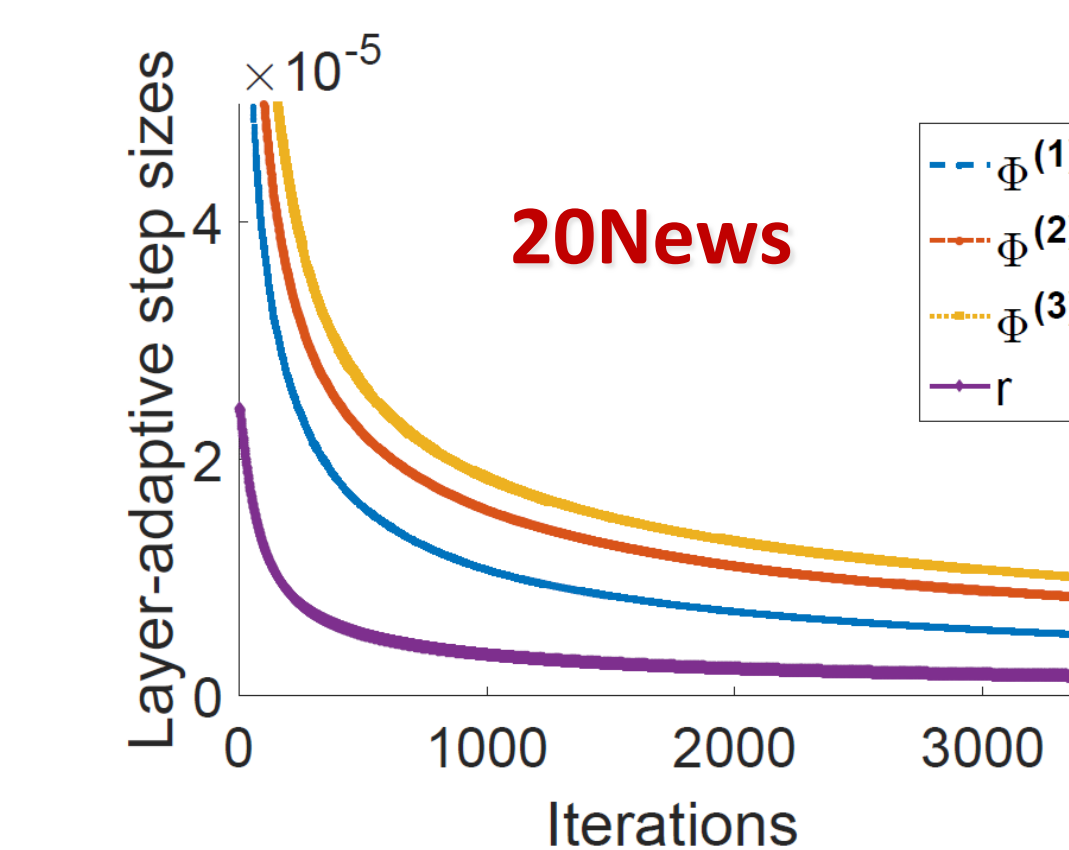


(b) DLDA of size 128-64 on RCV1



(c) DLDA of size 128-64 on Wiki

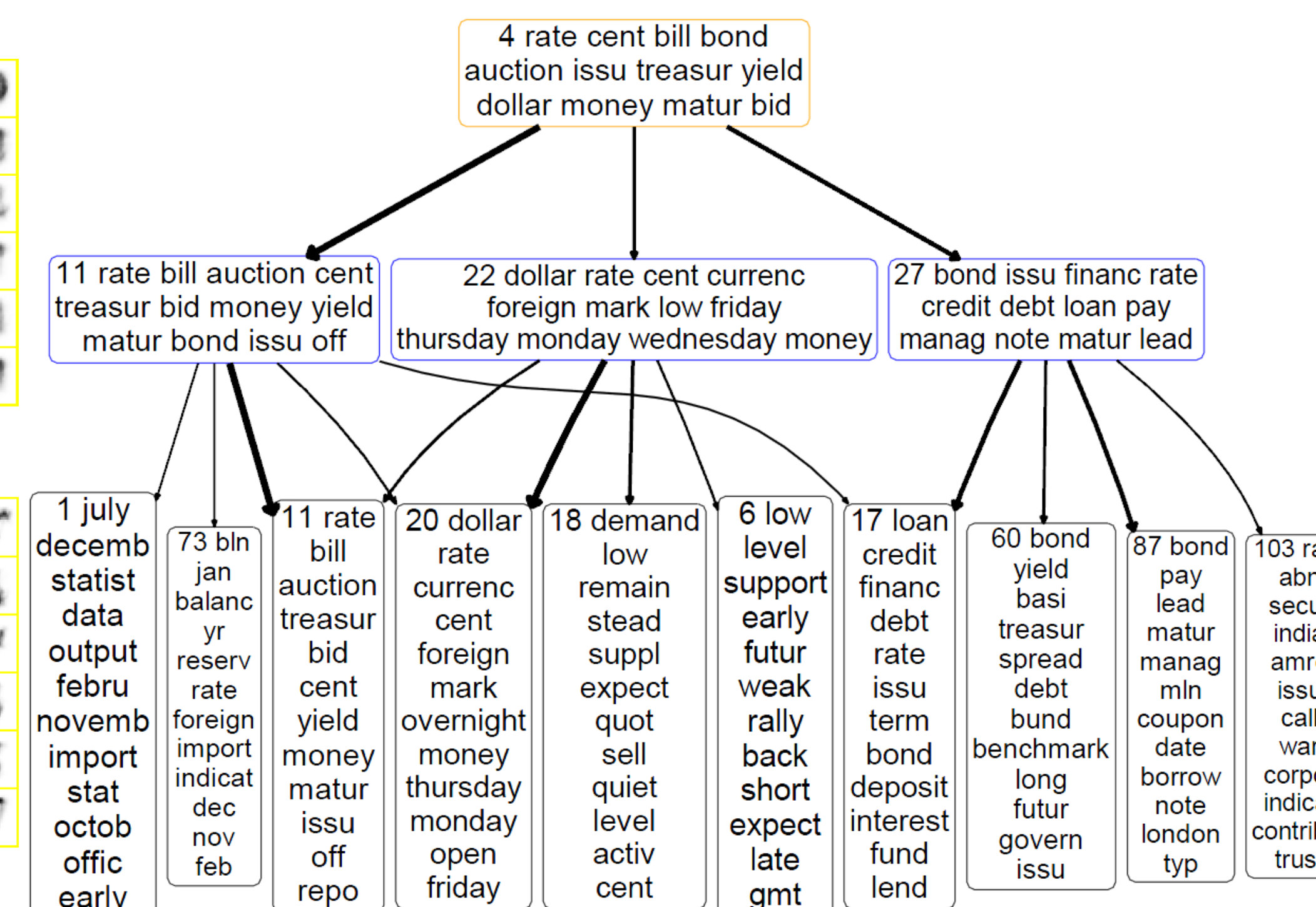
Scalable Joint Learning



Better Information Propagation



Example Topics for RCV1



Example Topics for Wiki

