

# Sparse Parallel Training of Hierarchical Dirichlet Process Topic Models

Alexander Terenin  
Imperial College London

Joint work with Måns Magnusson and Leif Jonsson

EMNLP 2020

November 16<sup>th</sup>–20<sup>th</sup>, 2020

[HTTPS://AVT.IM/](https://AVT.IM/) ·  @AVT\_IM



UPPSALA  
UNIVERSITET

Imperial College  
London

**A?**  
Aalto University

## Latent Dirichlet allocation

[DM Blei](#), [AY Ng](#), [MI Jordan](#) - [Journal of Machine Learning Research, 2003 - jmlr.org](#)

We describe latent Dirichlet allocation (LDA), a generative probabilistic model for collections of discrete data such as text corpora. LDA is a three-level hierarchical Bayesian model, in which each item of a collection is modeled as a finite mixture over an underlying set of ...

☆  Cited by 28439 [Related articles](#) [All 92 versions](#) 

## Hierarchical Dirichlet processes

[YW Teh](#), [MI Jordan](#), [MJ Beal](#), [DM Blei](#) - [Journal of the American Statistical ..., 2006 - JSTOR](#)

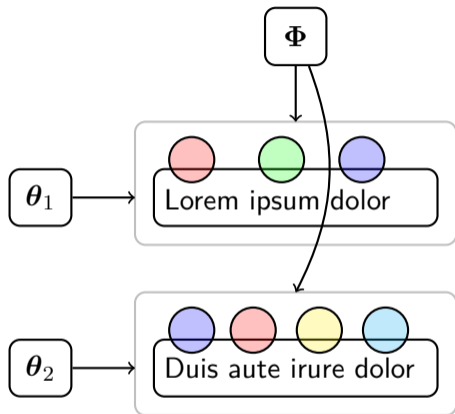
We propose the hierarchical Dirichlet process (HDP), a nonparametric Bayesian model for clustering problems involving multiple groups of data. Each group of data is modeled with a mixture, with the number of components being open-ended and inferred automatically by the ...

☆  Cited by 3853 [Related articles](#) [All 82 versions](#) 

The canonical topic models – everybody uses them!

This work: compute as fast, parallel, and principled as possible

# Latent Dirichlet Allocation



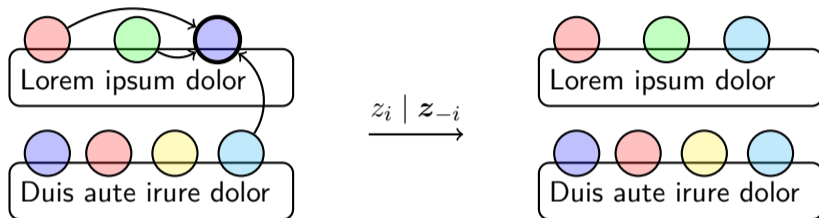
$\theta_d$  : document-topic probabilities

$z$  : topic indicators

$\Phi$  : topic-word probabilities

$w$  : word tokens

# Sparse Fully Collapsed Gibbs Sampling

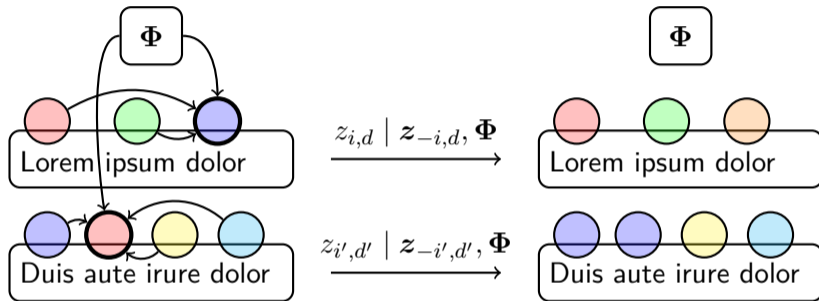


✓ Relatively fast

✗ Sequential

✗ Not fully sparse

# Sparse Partially Collapsed Gibbs Sampling



✓ Massively parallel

✗ Still not fully sparse

M. Magnusson, L. Jonsson, M. Villani, and D. Broman. Sparse Partially Collapsed MCMC for Parallel Inference in Topic Models. *Journal of Computational and Graphical Statistics* 27(2):449–463, 2018.

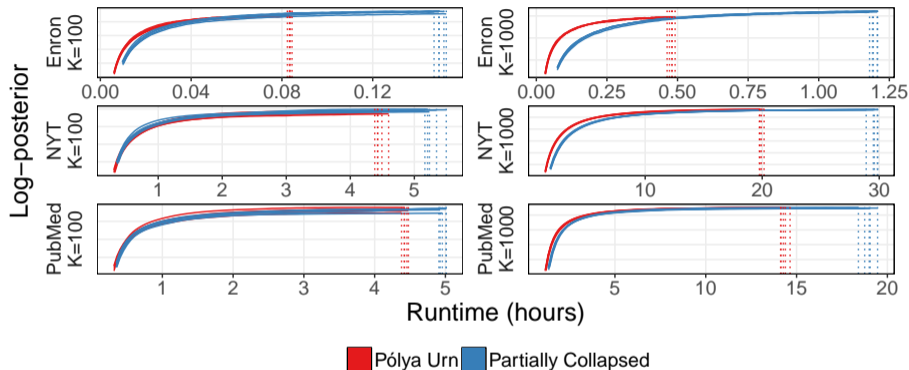
# Pólya Urn LDA

Dirichlet		Poisson Pólya Urn
$\mathbf{x} = \left[ \frac{\gamma_1}{\sum_{i=1}^k \gamma_i}, \dots, \frac{\gamma_k}{\sum_{i=1}^k \gamma_i} \right]$	$\rightarrow$	$\mathbf{y} = \left[ \frac{\tilde{\gamma}_1}{\sum_{i=1}^k \tilde{\gamma}_i}, \dots, \frac{\tilde{\gamma}_k}{\sum_{i=1}^k \tilde{\gamma}_i} \right]$
$\gamma_i \sim \text{Gamma}(\varpi F_i, 1)$		$\tilde{\gamma}_i \sim \text{Poisson}(\varpi F_i)$
Dense		Sparse

*Theorem.* Let  $\mathbf{x} \sim \text{Dir}(\varpi, \mathbf{F})$  and  $\mathbf{y} \sim \text{PPU}(\varpi, \mathbf{F})$ . Then for all  $\mathbf{F}$  we have  $d_{\text{LP}}(\mathbf{x}, \mathbf{y}) \rightarrow 0$  as  $\varpi \rightarrow \infty$ . ( $d_{\text{LP}}$ : Levy-Prokhorov metric)

- ✓ Dense matrix  $\Phi$  becomes sparse
- ✓ Utilize all sources of sparsity

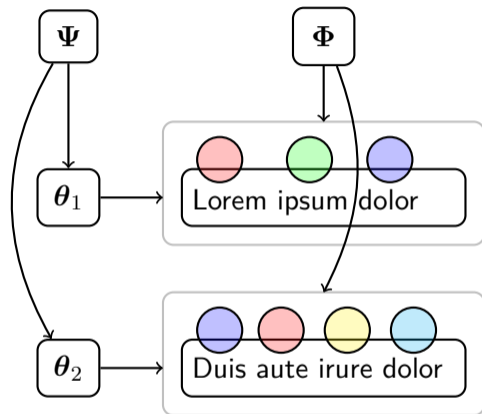
# Performance



✓ Faster runtime with no discernible loss in topic quality

A. Terenin, M. Magnusson, L. Jonsson, and D. Draper. Pólya Urn Latent Dirichlet Allocation: a doubly sparse massively parallel sampler. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41(7):1709–1719, 2019.

# Hierarchical Dirichlet Process Topic Model



$\theta_d$  : document-topic probabilities

$z$  : topic indicators

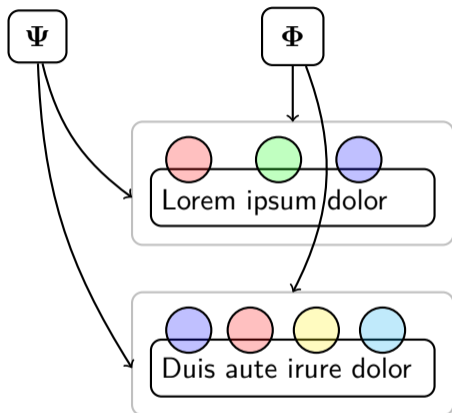
$\Phi$  : topic-word probabilities

$w$  : word tokens

$\Psi$  : global topic probabilities



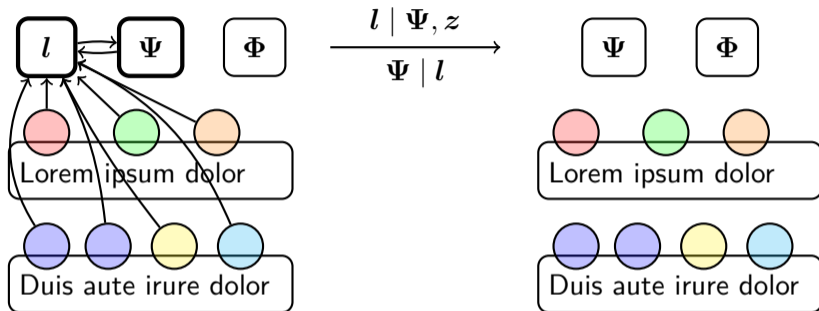
## Global topic probabilities



$\Psi$ : chosen to be a *Griffiths-Engen-McCloskey* (GEM) distribution

What's the posterior full conditional?

# Sparse Partially Collapsed Gibbs Sampling



*Theorem.*  $\Psi \mid l$  is given by

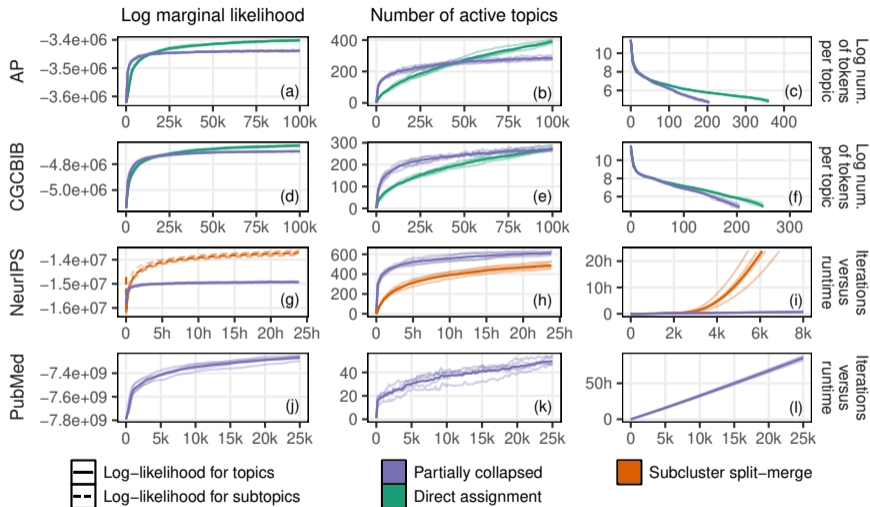
$$\Psi_k = s_k \prod_{i=1}^{k-1} (1 - s_i) \quad s_k \sim \text{B}(a_k^{(\Psi)}, b_k^{(\Psi)})$$

$$a_k^{(\Psi)} = 1 + l_k$$

$$b_k^{(\Psi)} = \gamma + \sum_{i=k+1}^{\infty} l_i$$

and each  $l_k$  is a sum of certain binomial random variables.

# Performance



✓ Sparsity and parallelism enable scalability to large data sets

# Topics for PubMed

## Large topics

47 322 709	40 229 486	34 685 122	30 795 166	30 707 144
care	age	model	cell	gene
health	risk	data	expression	protein
patient	children	system	growth	dna
medical	year	time	protein	expression
research	women	analysis	factor	sequence
clinical	patient	effect	receptor	genes
system	factor	test	kinase	rna
cost	population	field	beta	region

# Topics for PubMed

## Medium-size topics

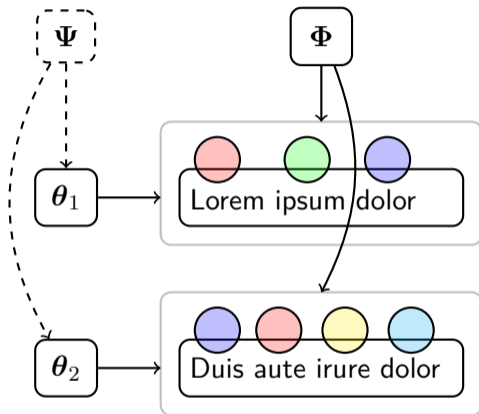
22 095 623	21 838 239	21 363 408	20 887 061	20 828 980
ca	structure	concentration	pregnancy	protein
effect	binding	degrees	level	binding
receptor	protein	samples	women	human
channel	reaction	liquid	hormone	antibodies
cell	acid	solution	day	acid
calcium	interaction	assay	fetal	alpha
concentration	compound	detection	infant	antibody
na	site	system	concentration	gel

# Topics for PubMed

## Small topics

16 136 164	14 201 063	13 706 016	13 191 158	13 105 245
effect	diet	patient	strain	protein
platelet	weight	disease	plant	membrane
induced	intake	gastric	growth	cell
oxide	food	asthma	acid	domain
rat	body	test	bacteria	binding
cell	effect	pylori	activity	receptor
endothelial	acid	arthritis	cell	lipid
activity	vitamin	chronic	species	membranes

## Prior information: HDP vs. LDA

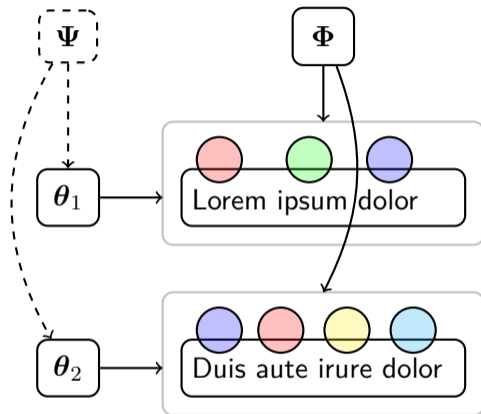


LDA:  $K$  topics with same number of tokens per topic *a priori*

HDP:  $\infty$  topics with number of tokens per topic decreasing

Number of active topics: proxy for true number of topics

## Prior information: HDP vs. LDA



Different prior specification yields different qualitative behavior

$$\Psi_k = s_k \prod_{i=1}^{\infty} (1 - s_k) \quad s_k \sim B(1, \gamma) \quad \text{vs.} \quad \Psi = U(1, \dots, K)$$



# Concluding remarks

Thank you for your attention!

[HTTPS://AVT.IM/](https://AVT.IM/) ·  @AVT\_IM

**Imperial College  
London**



UPPSALA  
UNIVERSITET



M. Magnusson, L. Jonsson, M. Villani, and D. Broman. Sparse Partially Collapsed MCMC for Parallel Inference in Topic Models. *Journal of Computational and Graphical Statistics* 27(2):449–463, 2018.

A. Terenin, M. Magnusson, L. Jonsson, and D. Draper. Pólya Urn Latent Dirichlet Allocation: a doubly sparse massively parallel sampler. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41(7):1709–1719, 2019.

A. Terenin, M. Magnusson, and L. Jonsson. Sparse Parallel Training of Hierarchical Dirichlet Process Topic Models. *Conference on Empirical Methods in Natural Language Processing*, 2020.