

# Leveraging Multi-Domain Prior Knowledge in Topic Models



Zhiyuan Chen<sup>†</sup> Arjun Mukherjee<sup>†</sup> Bing Liu<sup>†</sup>  
Meichun Hsu<sup>‡</sup> Malu Castellanos<sup>‡</sup> Riddhiman Ghosh<sup>‡</sup>

<sup>†</sup> University of Illinois at Chicago, <sup>‡</sup> HP Labs  
{czyuanacm, arjun4787}@gmail.com, liub@cs.uic.edu,  
meichun.hsu, malu.castellanos, riddhiman.ghosh}@hp.com

**UIC**  
UNIVERSITY  
OF ILLINOIS  
AT CHICAGO

## Introduction

- ❖ **Problem Definition:** Given prior knowledge from multiple domains, improve topic modeling in the **new** domain.
  - ❑ Knowledge in the form of **s-set** containing words sharing the same semantic meaning, e.g., {Light, Heavy, Weight}.
  - ❑ A novel technique to transfer knowledge to improve topic models.
- ❖ Existing Knowledge-based models
  - ❑ DF-LDA [Andrzejewski et al., 2009], Seeded Model (e.g., [Mukherjee and Liu, 2012]).
  - ❑ Two shortcomings: 1) Incapable of handling **multiple senses**, and 2) **Adverse effect** of Knowledge.

## MDK-LDA

- ❖ **Generative Process**
  - For each topic  $t \in \{1, \dots, T\}$ 
    - Draw a per topic distribution over s-sets,  $\varphi_t \sim \text{Dir}(\beta)$
    - For each s-set  $s \in \{1, \dots, S\}$ 
      - Draw a per topic, per s-set distribution over words,  $\eta_{t,s} \sim \text{Dir}(\gamma)$
  - For each document  $m \in \{1, \dots, M\}$ 
    - Draw  $\theta_m \sim \text{Dir}(\alpha)$
    - For each word  $w_{m,n}$ , where  $n \in \{1, \dots, N_m\}$ 
      - Draw a topic  $z_{m,n} \sim \text{Mult}(\theta_m)$
      - Draw an s-set  $s_{m,n} \sim \text{Mult}(\varphi_{z_{m,n}})$
      - Emit  $w_{m,n} \sim \text{Mult}(\eta_{z_{m,n}, s_{m,n}})$

### ❖ Plate Notation

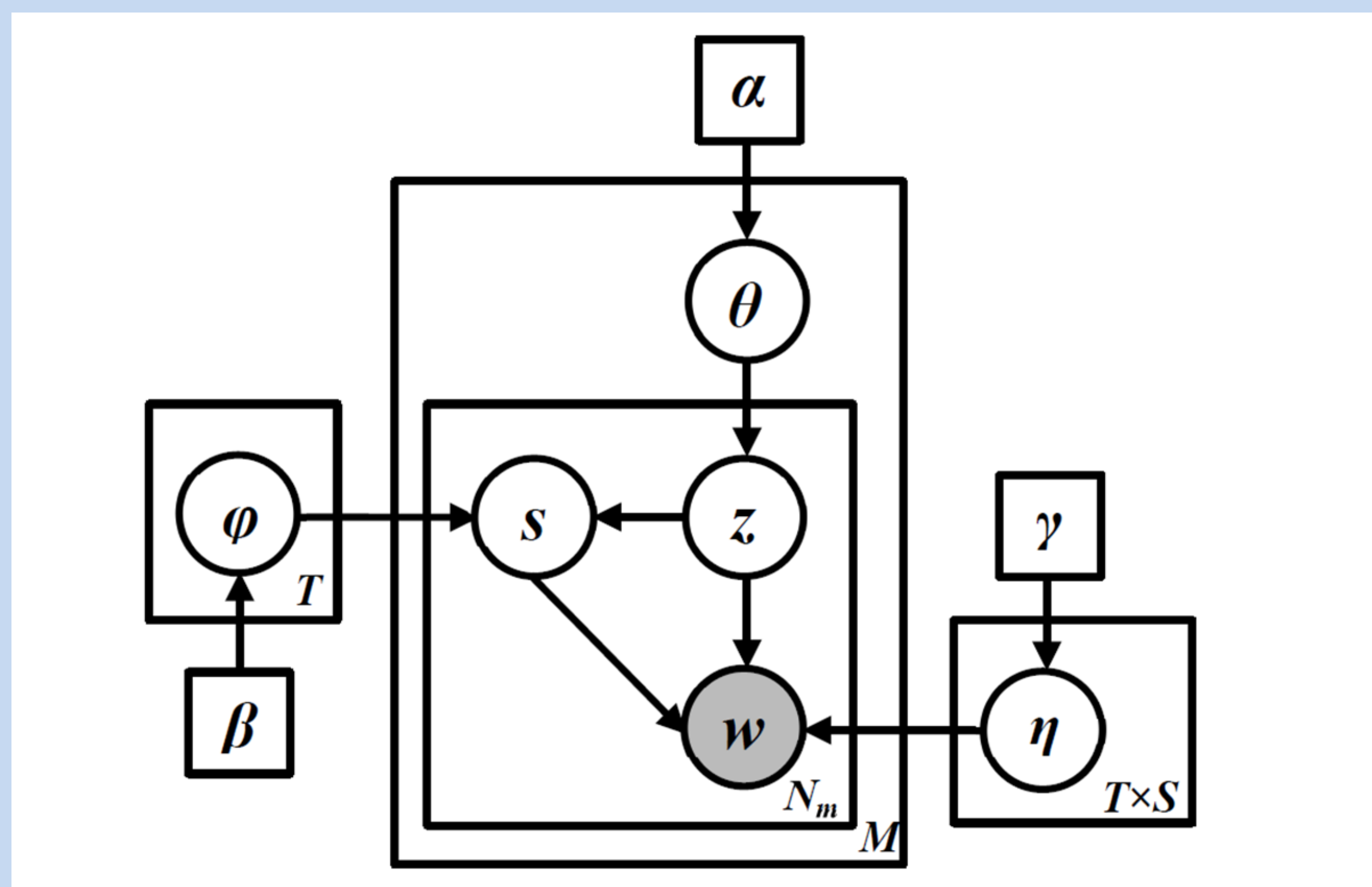


Figure 1: Plate notation of the proposed framework.

### ❖ Collapsed Gibbs Sampling

- ❑ Blocked Gibbs Sampler: Sample topic  $z$  and s-set  $s$  for word  $w$

$$P(z_i = t, s_i = s | \mathbf{z}^{-i}, \mathbf{s}^{-i}, \mathbf{w}, \alpha, \beta, \gamma) \propto \frac{n_{m,t}^{-i} + \alpha}{\sum_{t'=1}^T (n_{m,t'}^{-i} + \alpha)} \times \frac{n_{t,s}^{-i} + \beta}{\sum_{s'=1}^S (n_{t,s'}^{-i} + \beta)} \times \frac{n_{t,s,w_i}^{-i} + \gamma_s}{\sum_{v'=1}^V (n_{t,s,v'}^{-i} + \gamma_s)}$$

## Generalized Pólya Urn Model

- ❖ Generalized Pólya urn model [Mahmoud, 2008]
  - ❑ When a ball is drawn, that ball is put back along with a certain number of balls of **similar** colors.
- ❖ Promoting s-set as a whole
  - ❑ If a ball of color  $w$  is drawn, we put back  $\mathbb{A}_{s,w',w}$  balls of each color  $w' \in \{1, \dots, V\}$  where  $w$  and  $w'$  share s-set  $s$ .

$$\mathbb{A}_{s,w',w} = \begin{cases} 1 & w = w' \\ \sigma & w \in s, w' \in s, w \neq w' \\ 0 & \text{otherwise} \end{cases}$$

### ❖ Collapsed Gibbs Sampling

$$P(z_i = t, s_i = s | \mathbf{z}^{-i}, \mathbf{s}^{-i}, \mathbf{w}, \alpha, \beta, \gamma, \mathbb{A}) \propto \frac{n_{m,t}^{-i} + \alpha}{\sum_{t'=1}^T (n_{m,t'}^{-i} + \alpha)} \times \frac{\sum_{w'=1}^V \sum_{v'=1}^V \mathbb{A}_{s,v',w'} \cdot n_{t,s,v'}^{-i} + \beta}{\sum_{s'=1}^S (\sum_{w'=1}^V \sum_{v'=1}^V \mathbb{A}_{s',v',w'} \cdot n_{t,s',v'}^{-i} + \beta)} \times \frac{n_{t,s,w_i}^{-i} + \gamma_s}{\sum_{v'=1}^V (n_{t,s,v'}^{-i} + \gamma_s)}$$

## Experiments

- ❖ Datasets: reviews from six domains from Amazon.com.
- ❖ Baseline Models
  - ❑ **LDA** [Blei et al., 2003], **LDA\_GPU** [Mimno et al., 2011], and **DF-LDA** [Andrzejewski et al., 2009].
- ❖ Topic Discovery Results
  - ❑ Evaluation measure: **Precision @ n (p @ n)**.
  - ❑ Quantitative results in Table 1, Qualitative results in Table 2.
- ❖ Objective Evaluation
  - ❑ Topic Coherence [Mimno et al., 2011].

Domains	LDA	LDA GPU	DF-LDA	MDK-LDA(b)	MDK-LDA
Camera	0.80	0.50	0.67	0.81	<b>0.93</b>
Computer	0.67	0.60	0.56	0.70	<b>0.88</b>
Food	0.87	0.61	0.67	0.84	<b>0.91</b>
Care	0.81	0.64	0.72	0.92	<b>0.91</b>
Average	0.79	0.59	0.66	0.82	<b>0.91</b>

Table 1 (Quantitative): Avg. precision of each model across domains.

Camera (Battery)		Computer (Price)		Food (Taste)		Care (Tooth)	
LDA	MDK	LDA	MDK	LDA	MDK	LDA	MDK
battery	extra	<i>acer</i>	cheap	taste	flavor	<i>price</i>	tooth
<i>screen</i>	charge	<i>power</i>	price	salt	sweet	tooth	gum
life	life	<i>base</i>	inexpensive	<i>almond</i>	sugar	<i>amazon</i>	dentist
<i>lcd</i>	replacement	<i>year</i>	money	<i>fresh</i>	salty	pen	dental
<i>water</i>	battery	<i>button</i>	expensive	<i>pack</i>	tasty	<i>shipping</i>	whitening
usb	charger	<i>amazon</i>	cost	tasty	tasting	gum	pen
<i>cable</i>	<i>aa</i>	<i>control</i>	dollar	<i>oil</i>	delicious	dentist	refill
<i>case</i>	power	price	buck	<i>roasted</i>	taste	whitening	<i>year</i>
charger	rechargeable	<i>color</i>	worth	pepper	salt	refill	<i>date</i>
hour	time	purchase	low	<i>easy</i>	spice	<i>worth</i>	<i>product</i>

Table 2 (Qualitative): Example topics (MDK is short for MDK-LDA); **errors** are marked in red/italic.