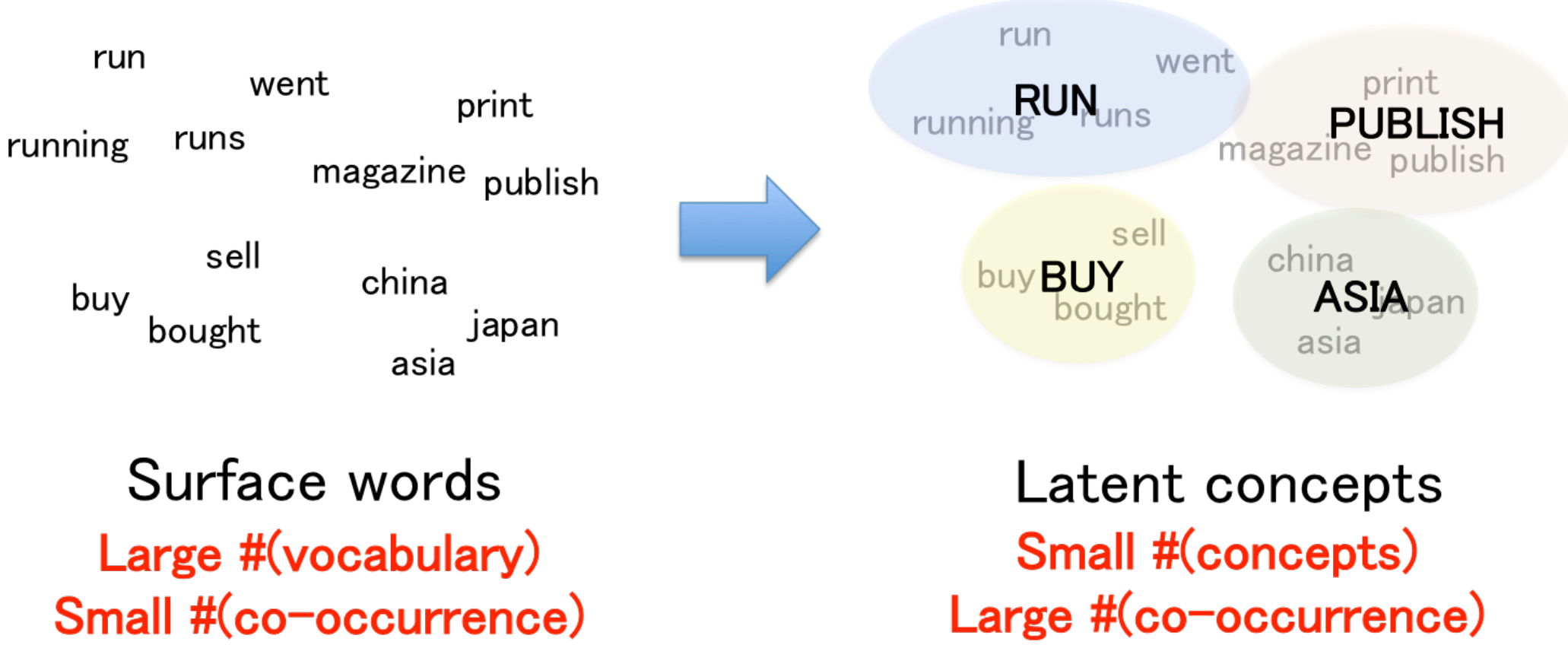


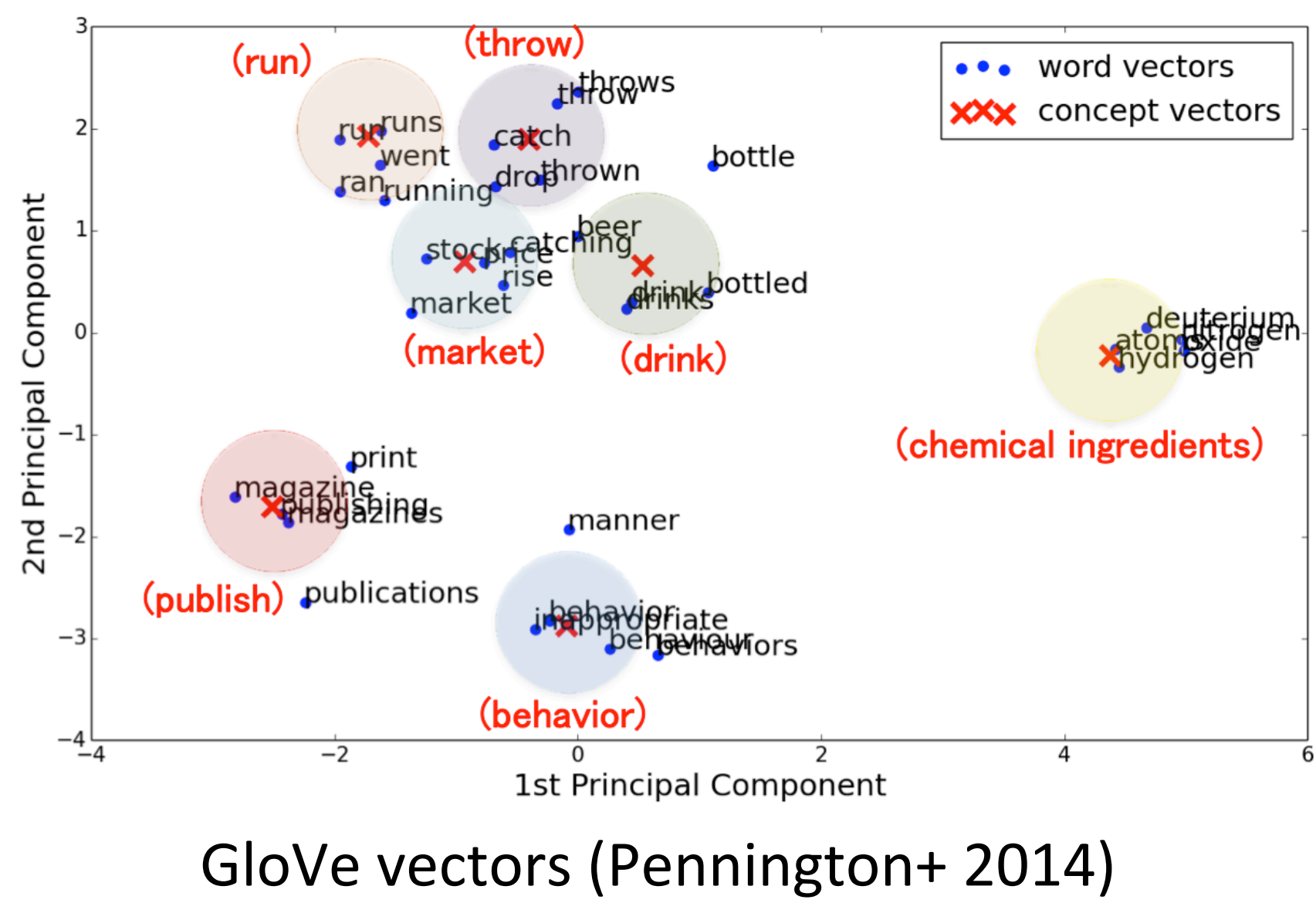
## Motivation

- Document-level word co-occurrence **is scarce when texts are short and vocabulary is diverse** (e.g. blog, SNS, newsgroup).
- Probabilistic topic models (e.g., LDA, pLSI) infers topics based on **document-level word co-occurrence**.
- Conventional topic models are not effective.
- Propose a novel topic model based on **co-occurrence statistics of latent concepts** to resolve the data sparsity.

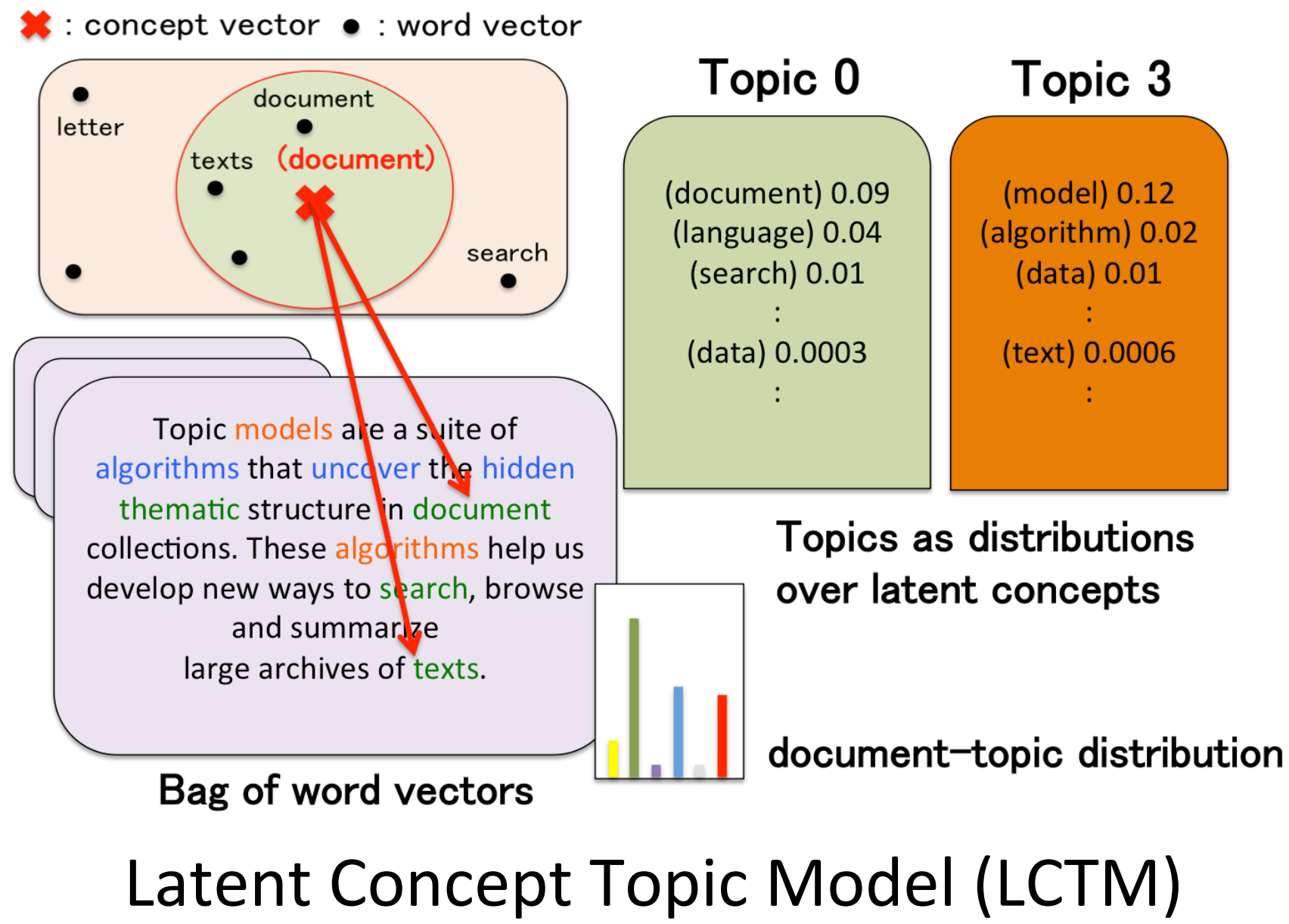


## Proposal

- Use **Neural word embedding** (e.g., word2vec, Glove) to capture **conceptual similarity of words**.
- Each cluster corresponds to one **latent concept**.



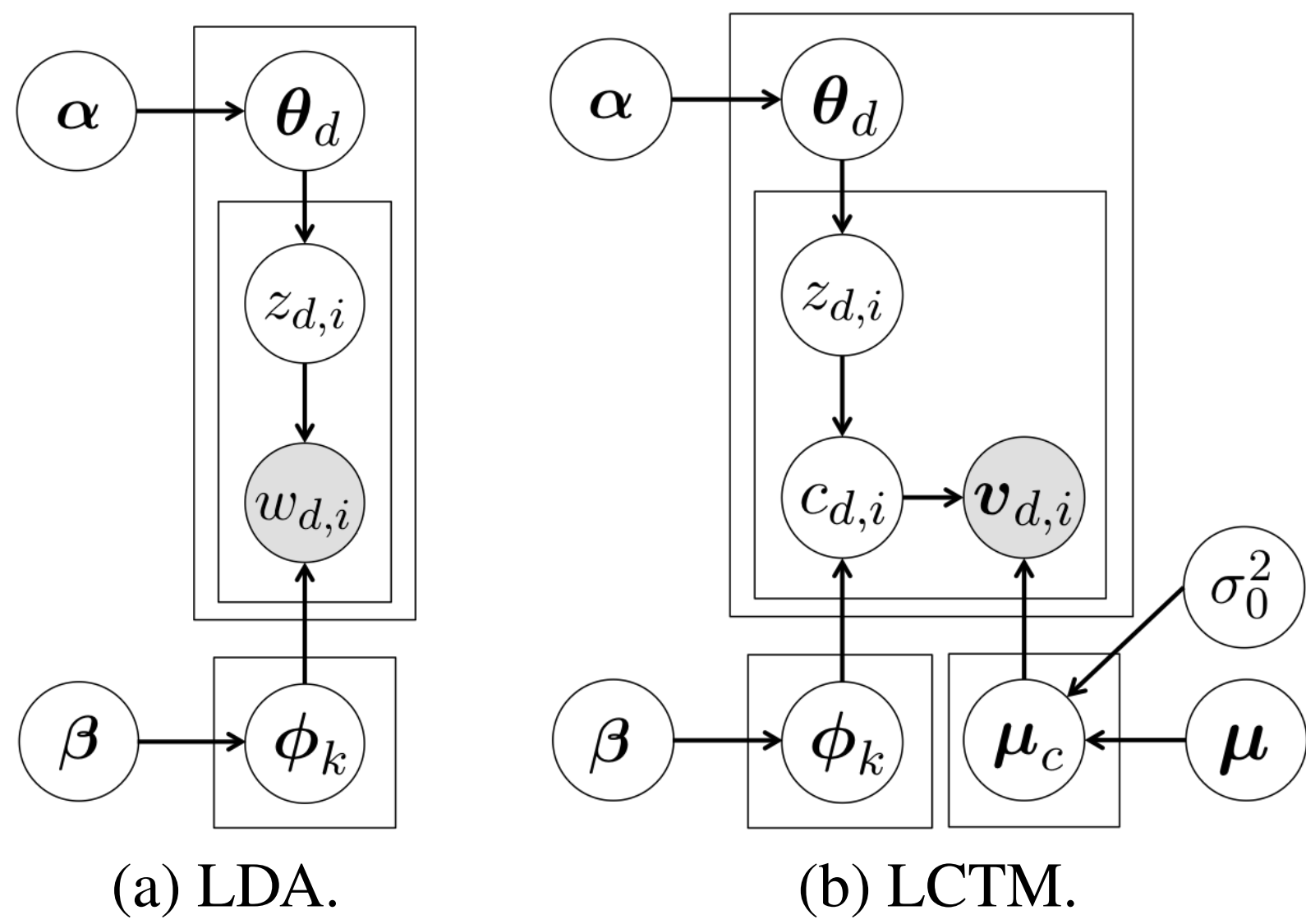
- Define **topics as distributions over latent concepts**.
- **Resolve data sparsity in short texts**.
- Model the **generative process of word embeddings**.
- LCTM can naturally handle **Out of Vocabulary (OOV) words**.



Gaussian variance parameter  $\sigma^2$  controls the range of the emission.

## Graphical Models

Add another layer of latent variables (**latent concepts**) to **mediate data sparsity**.



Notations

$\alpha$	$\beta$	Dirichlet prior parameters
$\mu$	$\sigma_0^2$	Gaussian prior parameters
$\theta_d$		document-topic distribution
$\phi_k$		topic-concept (word) distribution
$w_{d,i}$		word type
$v_{d,i}$		word vector
$z_{d,i}$		latent topic
$c_{d,i}$		latent concept
$\mu_c$		concept vector

## Overview of topic inference

- Collapsed Gibbs sampler** for the approximate inference.
- Sample **latent concepts** in addition to **topics**.

**Sampling of a topic assignment**

$$p(z_{d,i} = k \mid c_{d,i} = c, z^{-d,i}, c^{-d,i}, v) \propto \left( n_{d,k}^{-d,i} + \alpha_k \right) \cdot \frac{n_{k,c}^{-d,i} + \beta_c}{n_{k,\cdot}^{-d,i} + \sum_{c'} \beta_{c'}}$$

Prop of topic k in the same doc      Prob of topic k generating concept c

**Sampling of a concept assignment**

$$p(c_{d,i} = c \mid z_{d,i} = k, v_{d,i}, z^{-d,i}, c^{-d,i}, v^{-d,i}) \propto \left( n_{k,c}^{-d,i} + \beta_c \right) \cdot \mathcal{N}(v_{d,i} \mid \bar{\mu}_c, \sigma_c^2 I)$$

Prob of topic k generating concept c      Prob of concept c generating word vec v

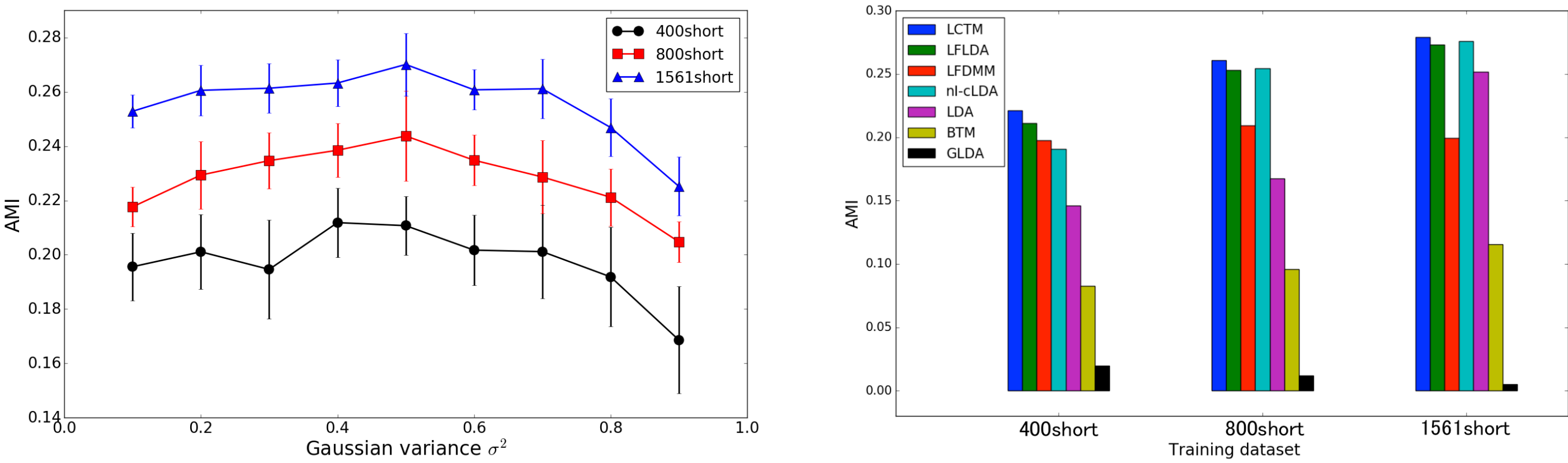
$\mathcal{N}(\cdot \mid \bar{\mu}_c, \sigma_c^2 I)$  : Gaussian distribution corresponding to latent concept c

## Experimental Results

**Dataset:** Short posts (less than 50 words) of 20Newsgroup.

### 1. Performance on document clustering

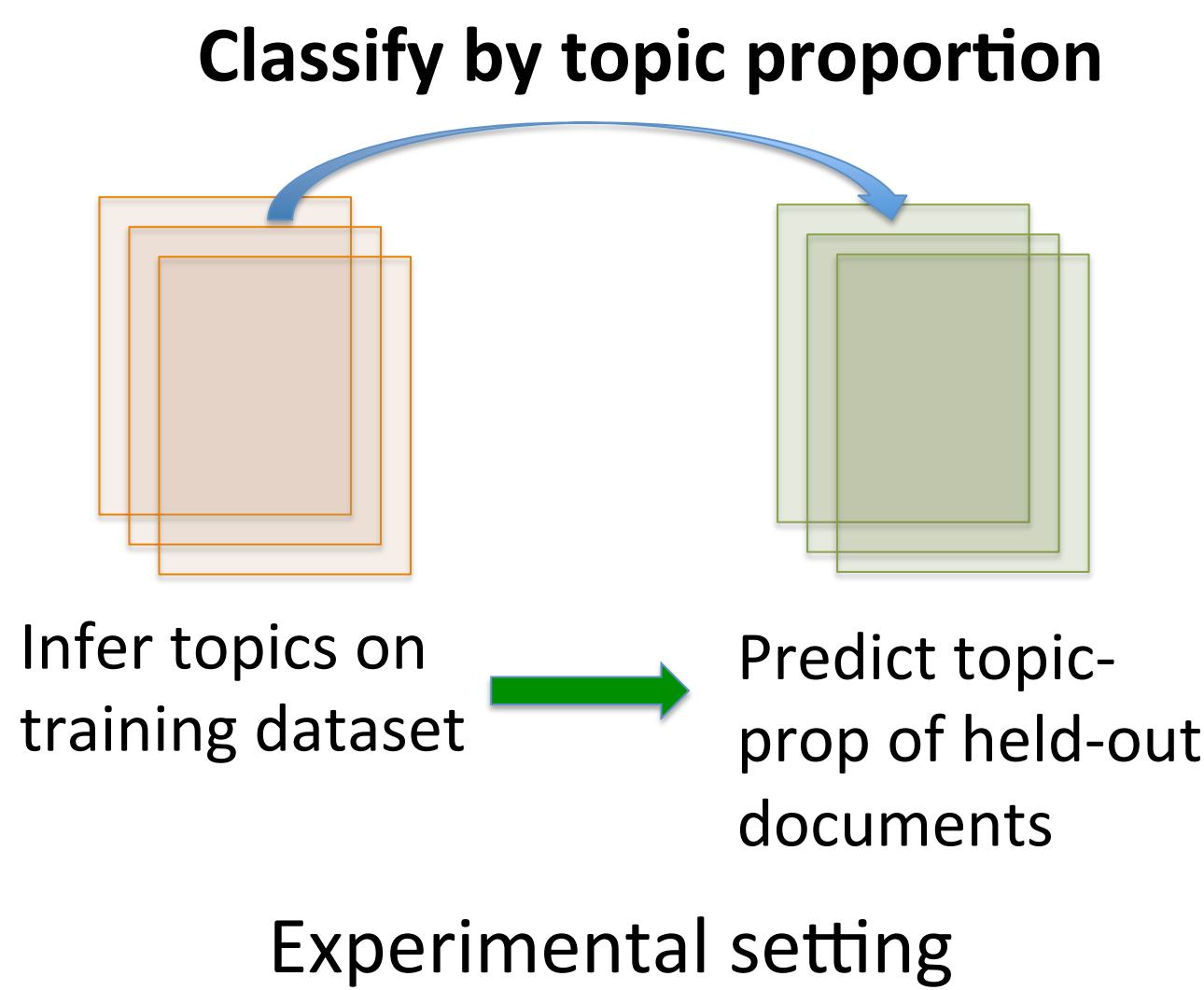
- Gaussian variance with  $\sigma^2 = 0.5$  consistently performs well.
- LCTM outperforms TM w/o word embeddings.
- LCTM performs comparable to TM w/ word embeddings.



Clustering performance measured by Adjusted Mutual Information (AMI)

### 2. Performance on handling a high degree of OOV words

- LCTM-UNK (LCTM that ignores OOV) outperforms other TMs.
- LCTM further improves performance of LCTM-UNK.
- LCTM effectively incorporates OOV words in held-out documents.



Training Set	400short	800short	1561short
OOV prop	0.348	0.253	0.181
Method	Classification Accuracy		
LCTM	<b>0.302</b>	<b>0.367</b>	<b>0.416</b>
LCTM-UNK	0.262	0.340	0.406
LFLDA	0.253	0.333	0.410
nI-cLDA	0.261	0.333	0.412
LDA	0.215	0.293	0.382
GLDA	0.0527	0.0529	0.0529
Chance Rate	0.0539	0.0539	0.0539

Classification accuracy on held-out documents

## Conclusion

- Introduced LCTM that infers topics based on document-level **co-occurrence of latent concepts**.
- Showed that LCTM can effectively **handle OOV words in held-out documents**.
- The same method can be readily applied to topic models that extend LDA.