



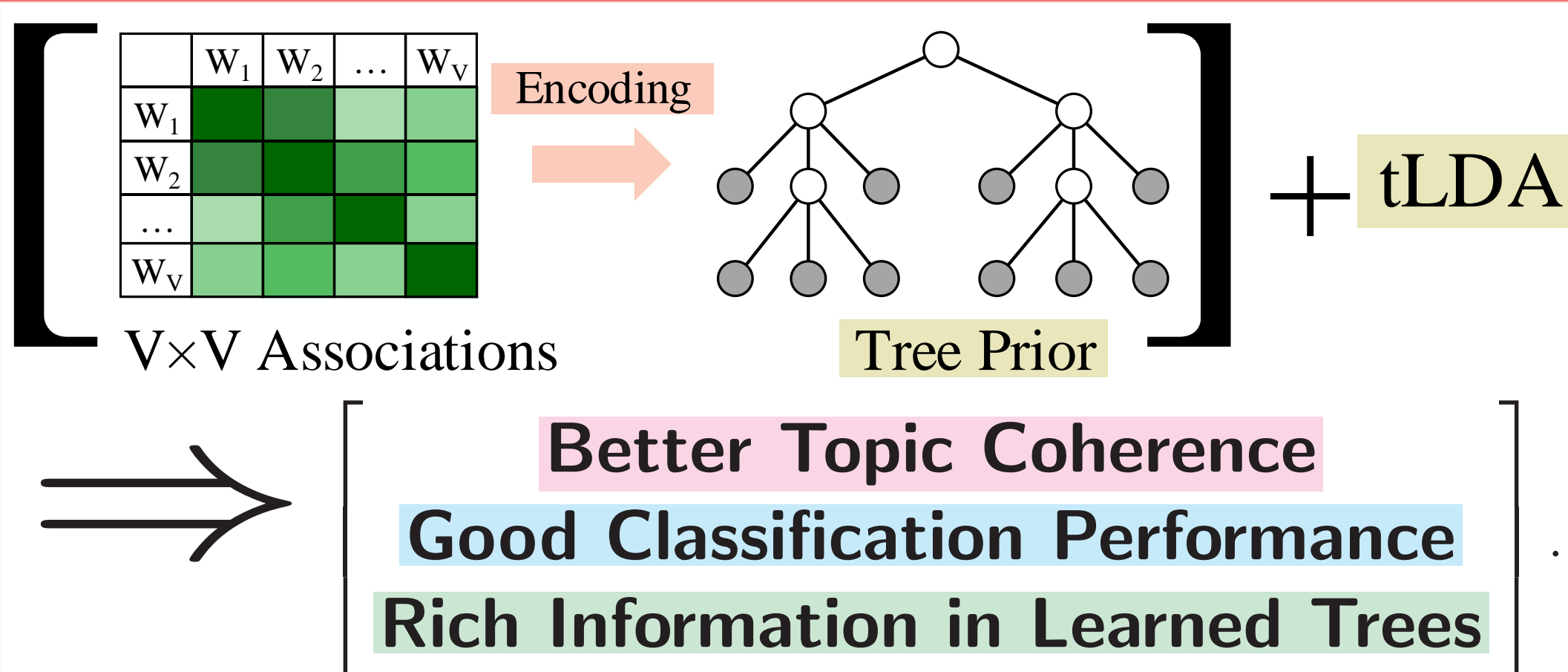
# Adapting Topic Models using Lexical Associations with Tree Priors

Weiwei Yang<sup>1,2</sup>, Jordan Boyd-Graber<sup>1,2,3,4</sup>, and Philip Resnik<sup>1,2,3,5</sup>

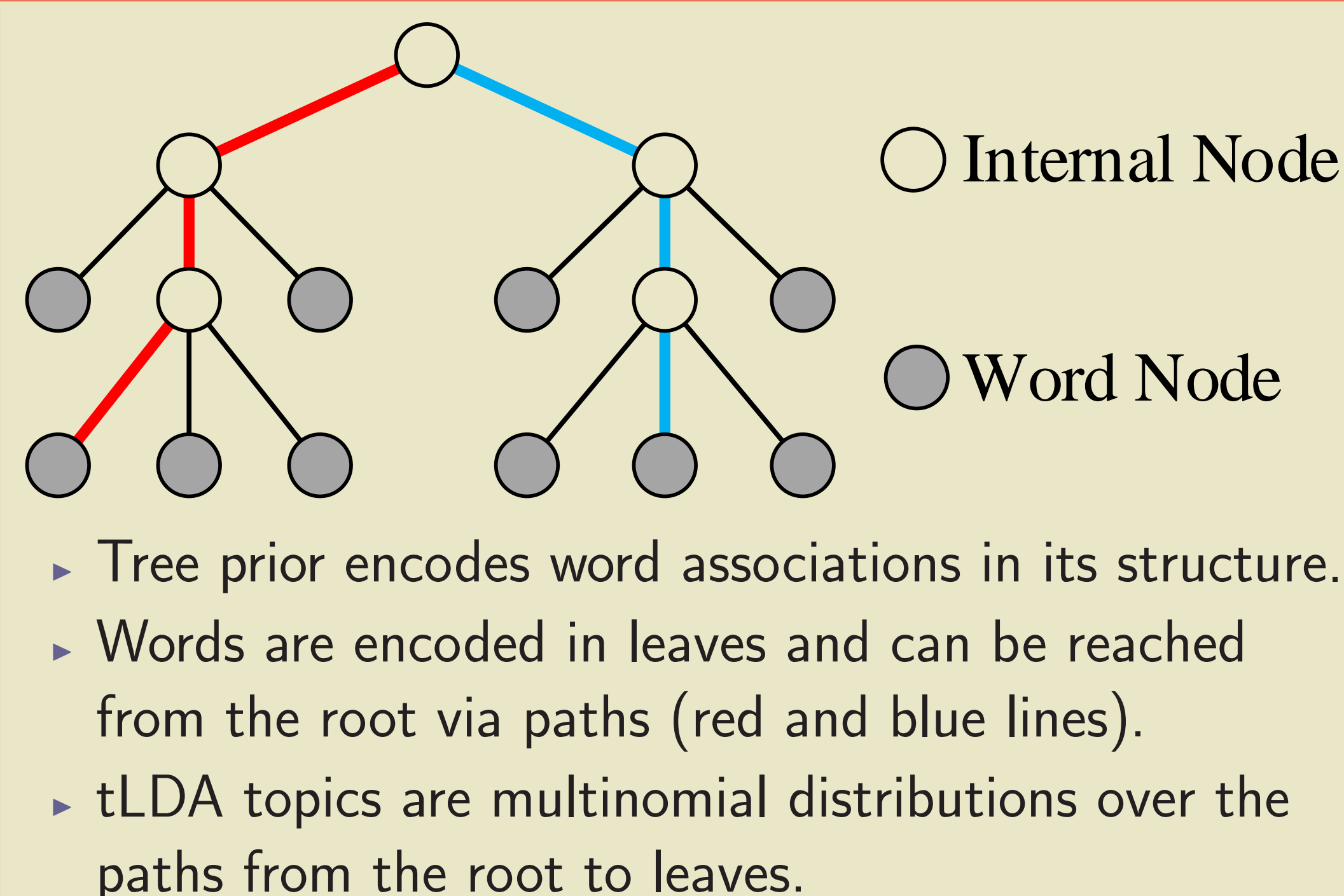
<sup>1</sup>Computer Science, <sup>2</sup>UMIACS, <sup>3</sup>Language Science Center, <sup>4</sup>iSchool, <sup>5</sup>Linguistics, University of Maryland, College Park



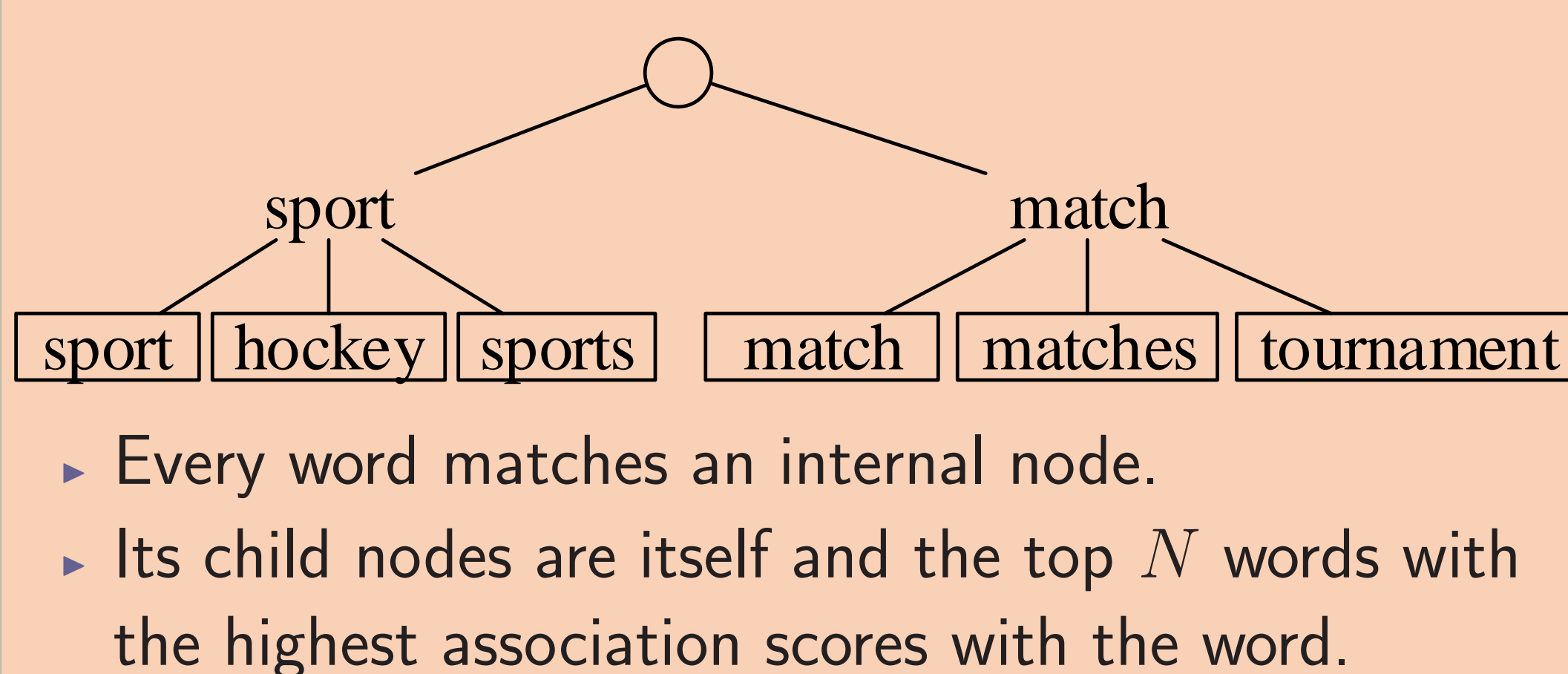
## Abstract



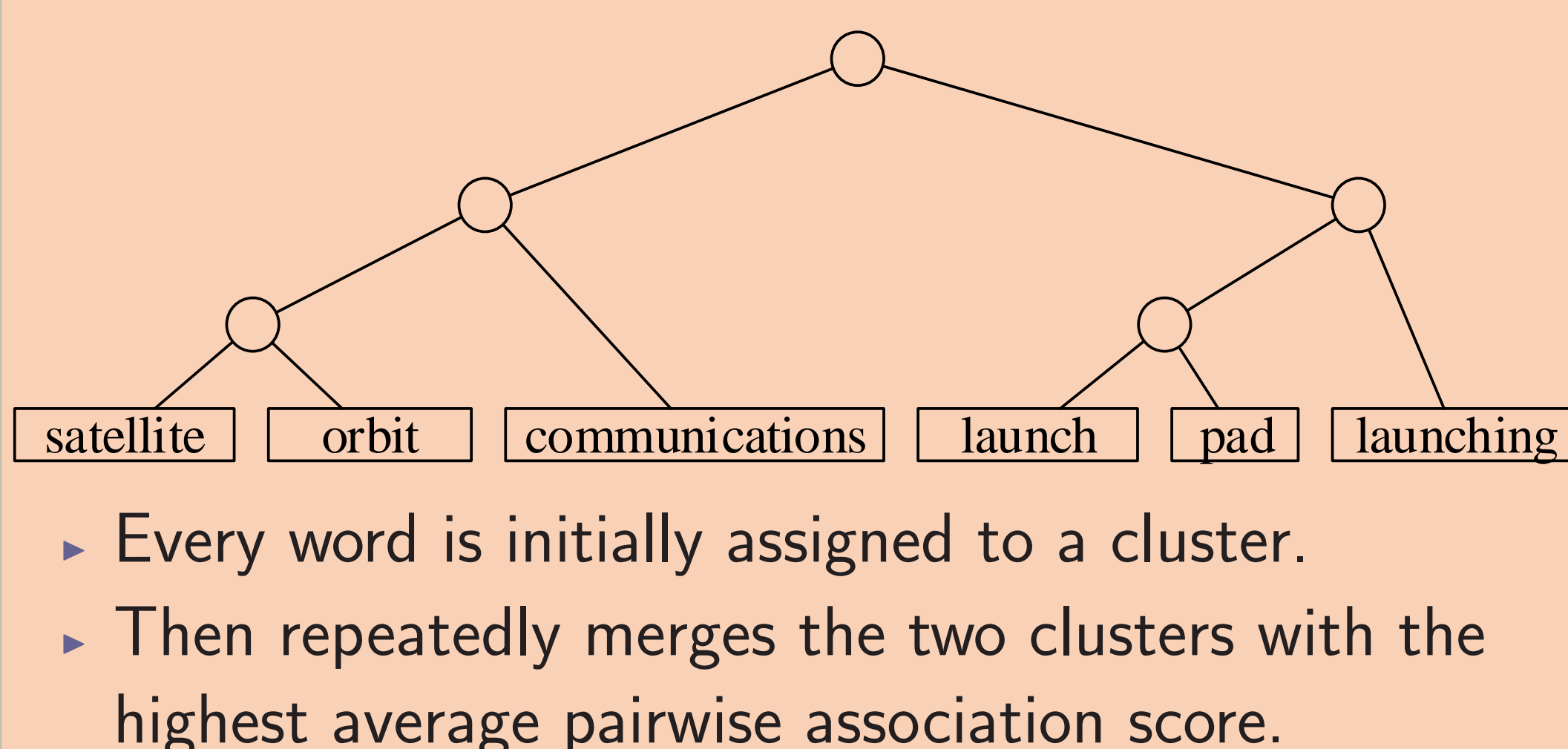
## Tree Prior and tLDA



## Two-Level (2LV)



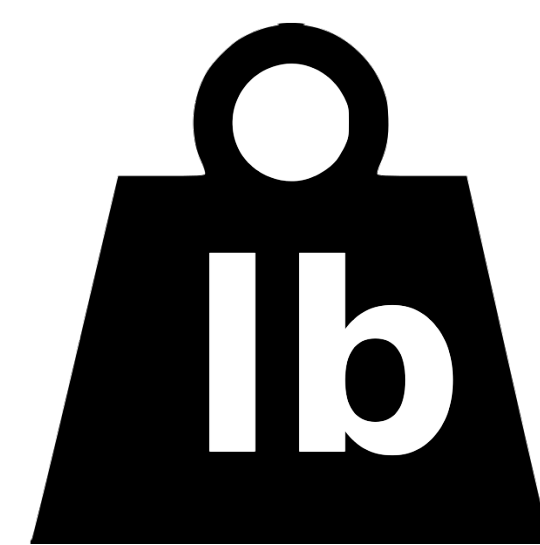
## Hierarchical Agglomerative Clustering (HAC)



## 1. Prepare Word Associations

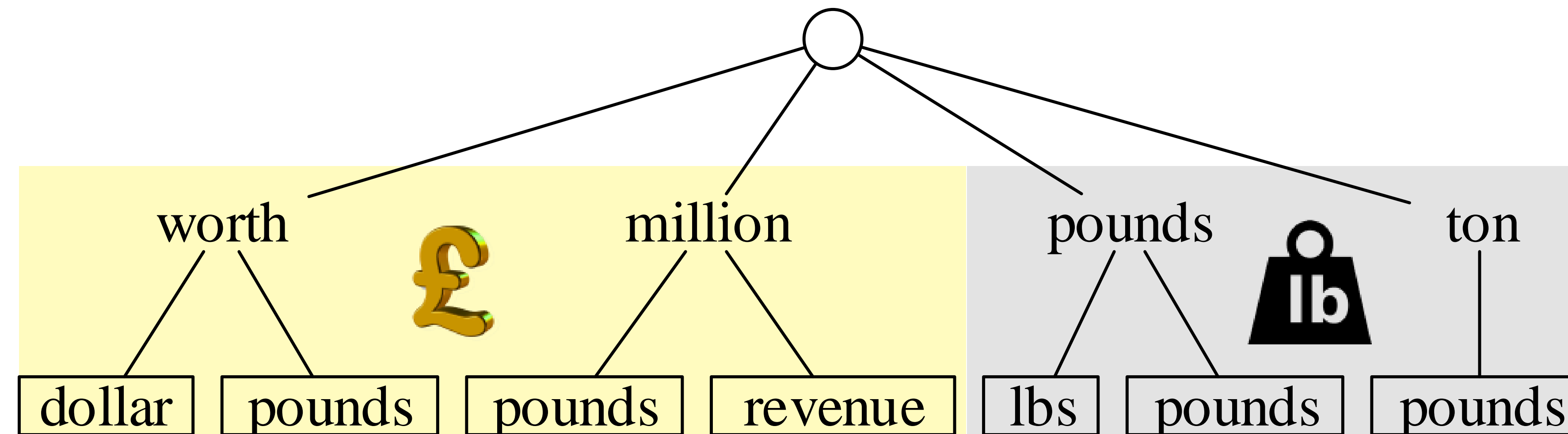
- Take "pounds" as an example.
- British currency vs. weight unit.

worth	dollar	.5529
worth	pounds	.5095
million	pounds	.5161
million	revenue	.5628
lbs	pounds	.6070
ton	pounds	.4767

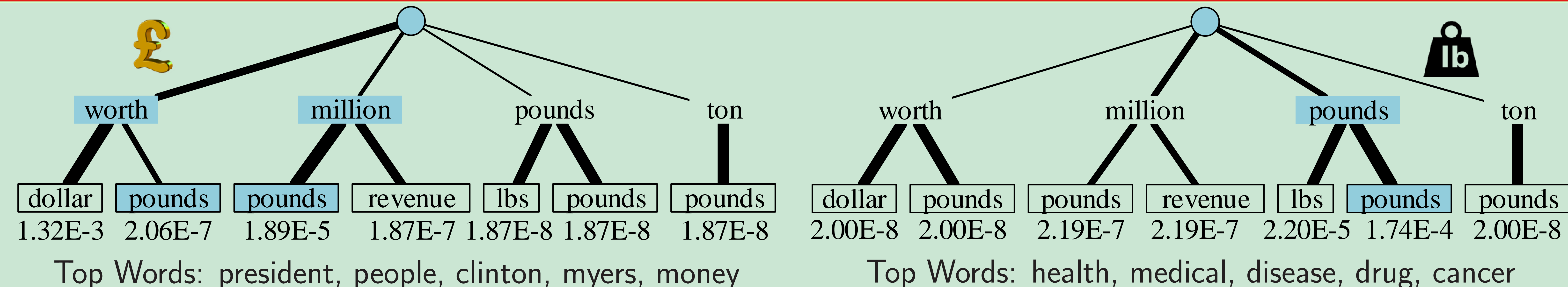


## Pipeline

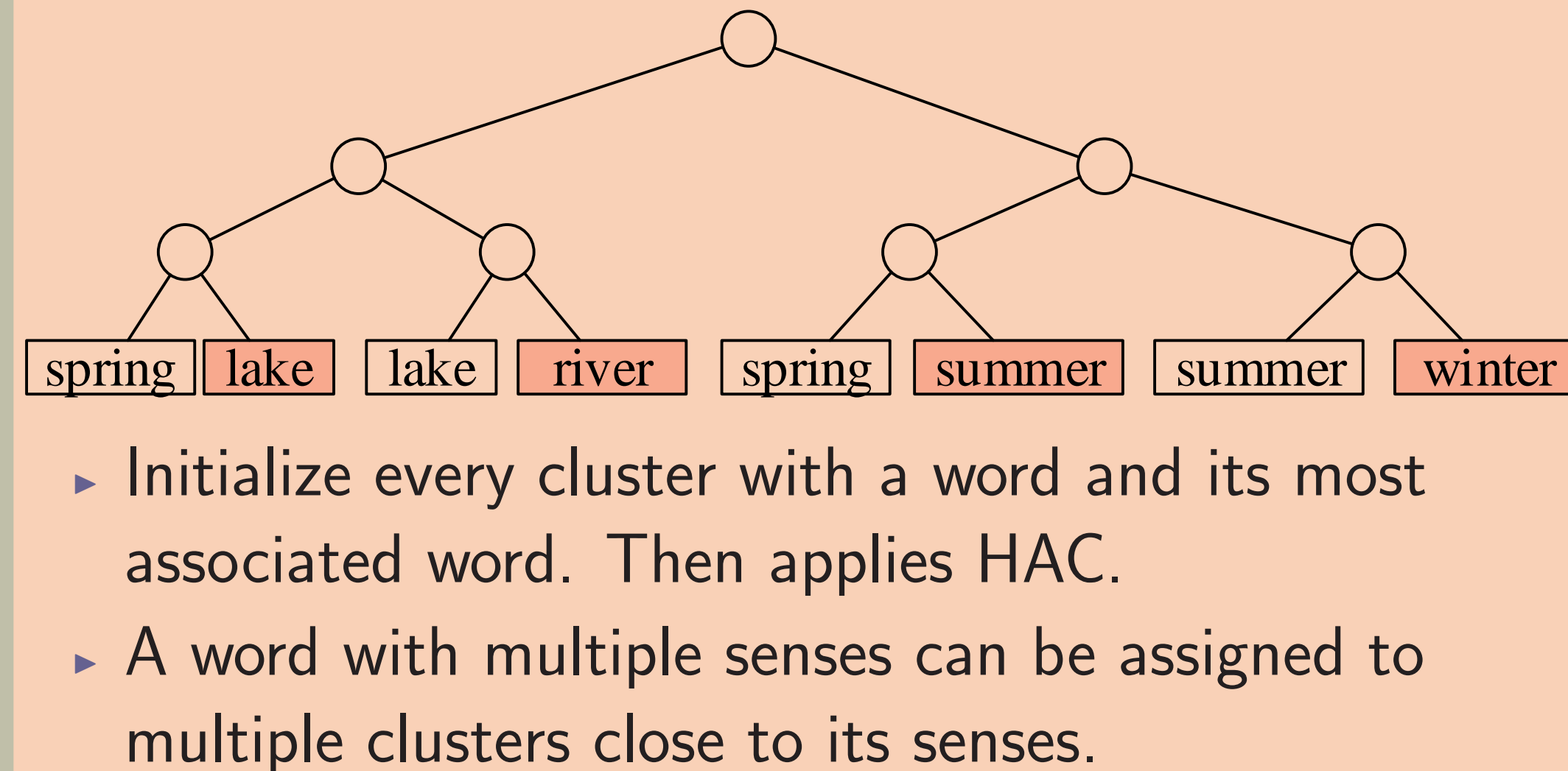
### 2. Construct Tree Priors



### 3. Learn Topics



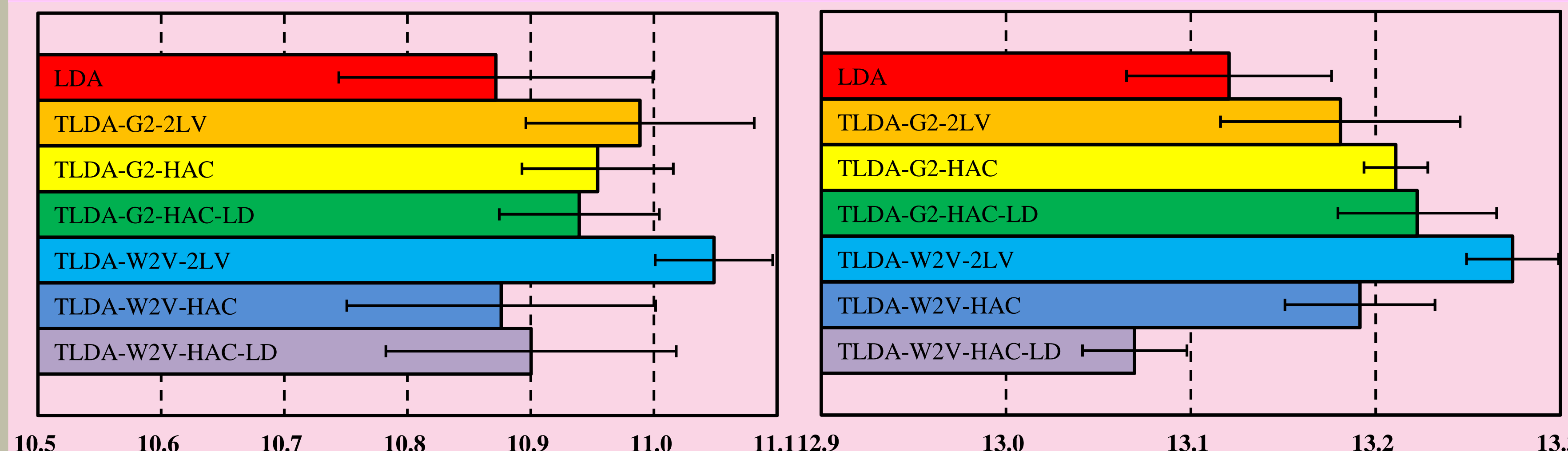
## HAC with Leaf Duplication (HAC-LD)



## Topic Coherence

- Datasets:** 20NewsGroups (20NG, left) and Amazon product reviews (right).
- Associations:** word2vec (W2V) and Dunning likelihood (G2).
- Metric:** Average pairwise PMI value of models' topics' top 10 words, on Wikipedia reference corpus. Higher is better.
- Baselines:** LDA and latent concept topic model (LCTM).
- Summary:** tLDA generally yields more coherent topics quantitatively. LCTM performs too poorly to be included.
- For topic words and qualitative analysis, see the paper. ↓

## Topic Coherence



## Paper



## Classification

Model	Tree	Path	20NG	Amazon
BOW	–	–	86.64	86.73
BOW+VEC	–	–	86.59	<b>87.30</b>
LDA	–	–	86.67	86.99
LCTM	–	–	86.52	86.83
tLDA	2LV	N	86.75	87.07
(W2V)	HAC	–	86.79	87.19
	HAC-LD	N	86.73	87.02
		Y	86.94	86.88
tLDA	2LV	N	86.82	87.15
(G2)		Y	<b>86.96</b>	87.05
	HAC	–	86.63	87.11
	HAC-LD	N	86.73	87.07
		Y	86.91	86.94

- 20NG:** Multi-class SVM classification. Documents' groups are their labels.
- Amazon:** Binary SVM classification. 4-5 stars are positive and 1-2 stars are negative.