

# Scale-Invariant Infinite Hierarchical Topic Model

Shusei Eshima

Department of Government, Harvard University  
shuseieshima@g.harvard.edu

Daichi Mochihashi

The Institute of Statistical Mathematics  
daichi@ism.ac.jp

## Problem Statement

- Topic models: summarize, annotate, and categorize documents for a human reader
- Hierarchical topic models: learning hierarchical latent topic organization, especially for a large number of topics (say, > 1000) [1, 2, 3]

Existing models: **fragmented tree structures**

- Expected probabilities of topics decay exponentially along the depth of tree
- Heuristic rules to update topics
- Restricting the tree structure (e.g., truncating the depth to three levels)

## Contributions

- Adjusting the probability scale by considering the size of the parent topic → less fragmentation
- Extending the tree-structured stick-breaking process (TSSB) [4] → variety of applications
- Efficiently drawing an infinite topic tree for each document from a base infinite tree in a hierarchical Bayesian fashion

## Scale-Invariant TSSB

TSSB is a crucial building block of recent neural hierarchical topic models [2, 3]. Let  $\kappa < \epsilon$  indicate that  $\kappa$  is an ancestor of  $\epsilon$ ,

$$\pi_\epsilon = \nu_\epsilon \prod_{\kappa < \epsilon} (1 - \nu_\kappa) \cdot \prod_{\kappa \leq \epsilon} \phi_\kappa, \quad \phi_{\epsilon\kappa} = \psi_{\epsilon\kappa} \prod_{j=1}^{k-1} (1 - \psi_{\epsilon_j})$$

$$\nu_\epsilon \sim \text{Be}(1, \alpha_0), \quad \psi_\epsilon \sim \text{Be}(1, \gamma_0).$$

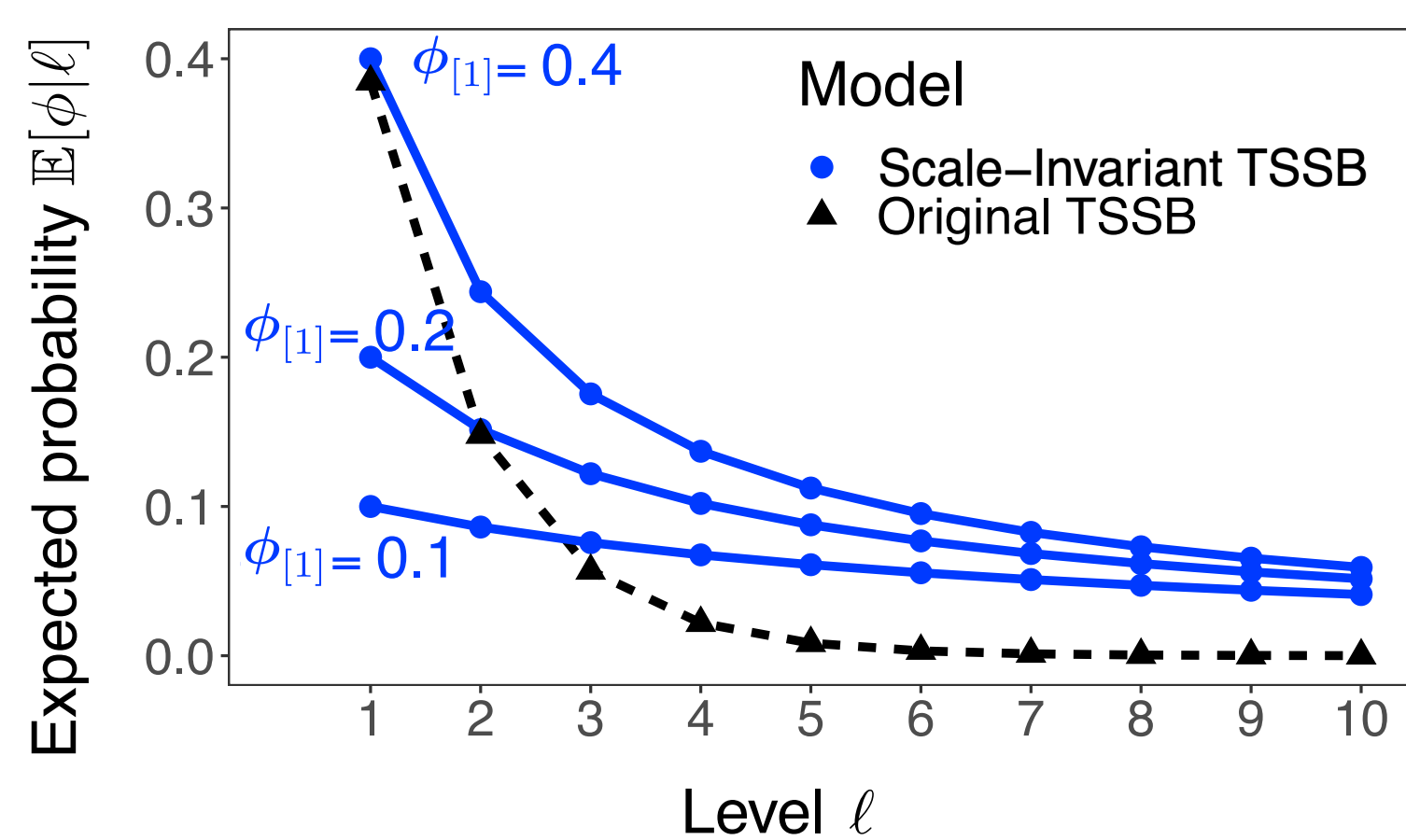
- The first term: the probability of stopping at the topic  $\epsilon$  vertically.
- The next product terms: to passing ancestors of  $\epsilon$  while horizontally stopping at  $\epsilon$  and its ancestors.
- Vertical and horizontal probabilities of stopping follow Beta distribution.

The scale-invariant TSSB rescales  $\gamma_0$ ,

$$\psi_\epsilon \sim \text{Be}(1, \phi_{\epsilon'} \gamma_0).$$

Using the horizontal breaking proportion of a parent topic,  $\phi_{\epsilon'}$ , to draw a *relative* stick length for its child topic → A larger break if the stick to break is shorter. The expected probability of horizontal break at  $\ell$  is  $\mathbb{E}[\phi|\ell] \approx 1/(2\gamma+1/\mathbb{E}[\phi|\ell-1])$  for  $\ell \geq 2$ . It was originally  $\mathbb{E}[\phi|\ell] \approx 1/(2\gamma+1)^\ell$ .

SI-TSSB has a slower decay.



## References

- D. M. Blei, T. L. Griffiths, M. I. Jordan, and J. B. Tenenbaum, "Hierarchical topic models and the nested Chinese restaurant process," in *NIPS*, 2003, pp. 17–24.
- M. Isonuma, J. Mori, D. Bollegala, and I. Sakata, "Tree-structured neural topic model," in *ACL*, 2020, pp. 800–806.
- Z. Chen, C. Ding, Z. Zhang, Y. Rao, and H. Xie, "Tree-structured topic modeling with nonparametric neural variational inference," in *ACL*, 2021, pp. 2343–2353.
- R. P. Adams, Z. Ghahramani, and M. I. Jordan, "Tree-structured stick breaking for hierarchical data," in *NIPS*, 2010, pp. 19–27.
- Y. W. Teh, "A Bayesian interpretation of interpolated Kneser-Ney," NUS School of Computing, Tech. Rep., 2006.

## Scale-Invariant Infinite Hierarchical Topic Model (ihLDA)

### Document-Topic Distribution:

- We apply the hierarchical Dirichlet process separately to the vertical and horizontal probabilities.

$$\nu_\epsilon \sim \text{Be}(a\tau_{\bar{\epsilon}}, a(1 - \sum_{\kappa \leq \bar{\epsilon}} \tau_\kappa)) \quad \text{where } \tau_\epsilon = \nu_\epsilon \prod_{\kappa < \epsilon} (1 - \nu_\kappa)$$

$$\psi_{\epsilon\kappa} \sim \text{Be}(b\phi_{\bar{\epsilon}\kappa}, b(1 - \sum_{j=1}^k \phi_{\bar{\epsilon}j}))$$

### Topic-Word Distribution:

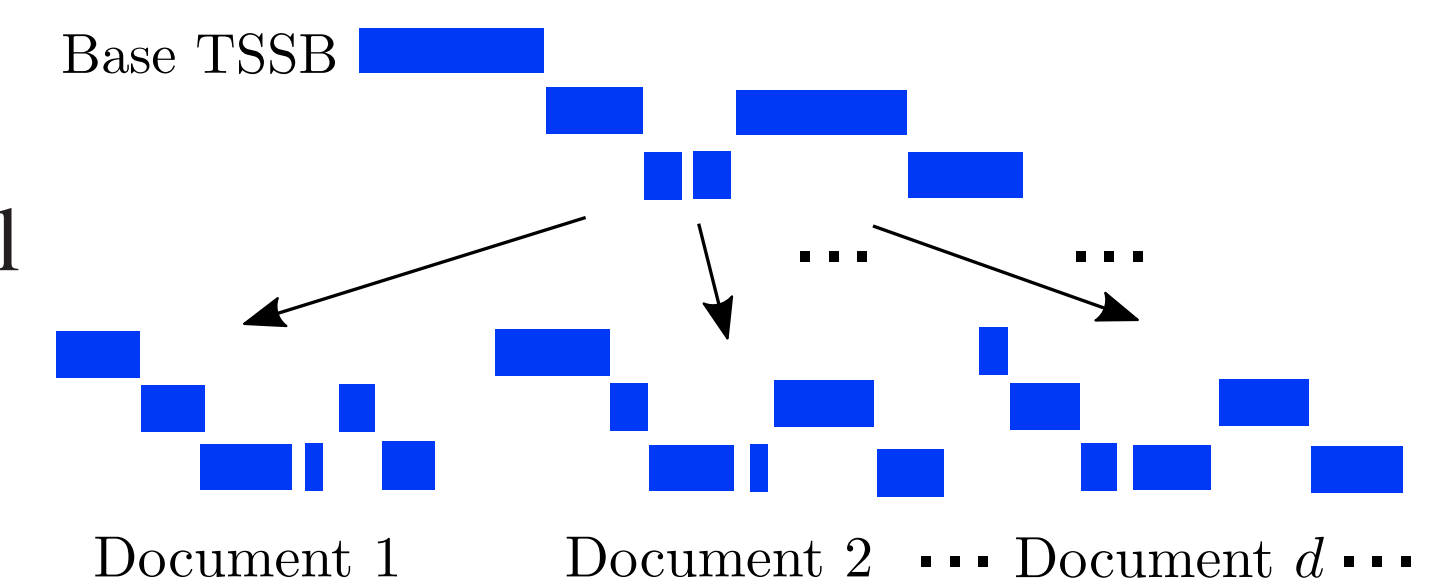
- The hierarchical Pitman-Yor process [5]
- The semantic similarity between a parent topic and its children
- Increasing the specificity as the tree deepens

### Inference:

- Gibbs sampling, Retrospective sampling + Binary search, Slice sampling

### Data Generation Process: $\pi^{(d)}$ := a TSSB for a doc $d$ .

- Draw a base TSSB  $\tilde{\pi}$ .
- Draw topic-word distributions  $H_\epsilon$  from the HPY for each topic in  $\tilde{\pi}$ .
- Draw a document-topic distribution for each document  $d$ ,  $\pi^{(d)} \sim \text{HTSSB}(\tilde{\pi})$ .
- For each word position  $i$  in a document  $d$ , draw a topic,  $z_{di} \sim \pi^{(d)}$  and draw a word,  $w_{di} \sim H_{z_{di}}$ .



## Comparing Top Words from Three Models (max level = 3 for comparison)

Proposed:	Probabilistic: nCRP [1]	Neural: TSNTM [2]
L1: said year would also people	L1: said year one time would	L1: said show year also would
L2: said people mobile technology phone	L2: said year also would company	L2: said year game world time
L3: said software site user mail	L3: film show magic would child	L3: england first game ireland win
L2: said would government people law	L3: film indian star india actor	L3: said labour blair party election
L3: tax said government would budget	L3: film dvd effect extra man	L3: said would people law gov.
L3: labour election said party blair	L3: film harry potter dvd warner	L3: said would gov. election tax
L2: film said best award year	L2: best award film actor actress	L3: said would tax gov. election
L3: music band song year album	L2: film star story life singer	L3: said would tax gov. election
L3: game dvd film year sony	L2: film star movie actress also	L3: said would tax gov. election

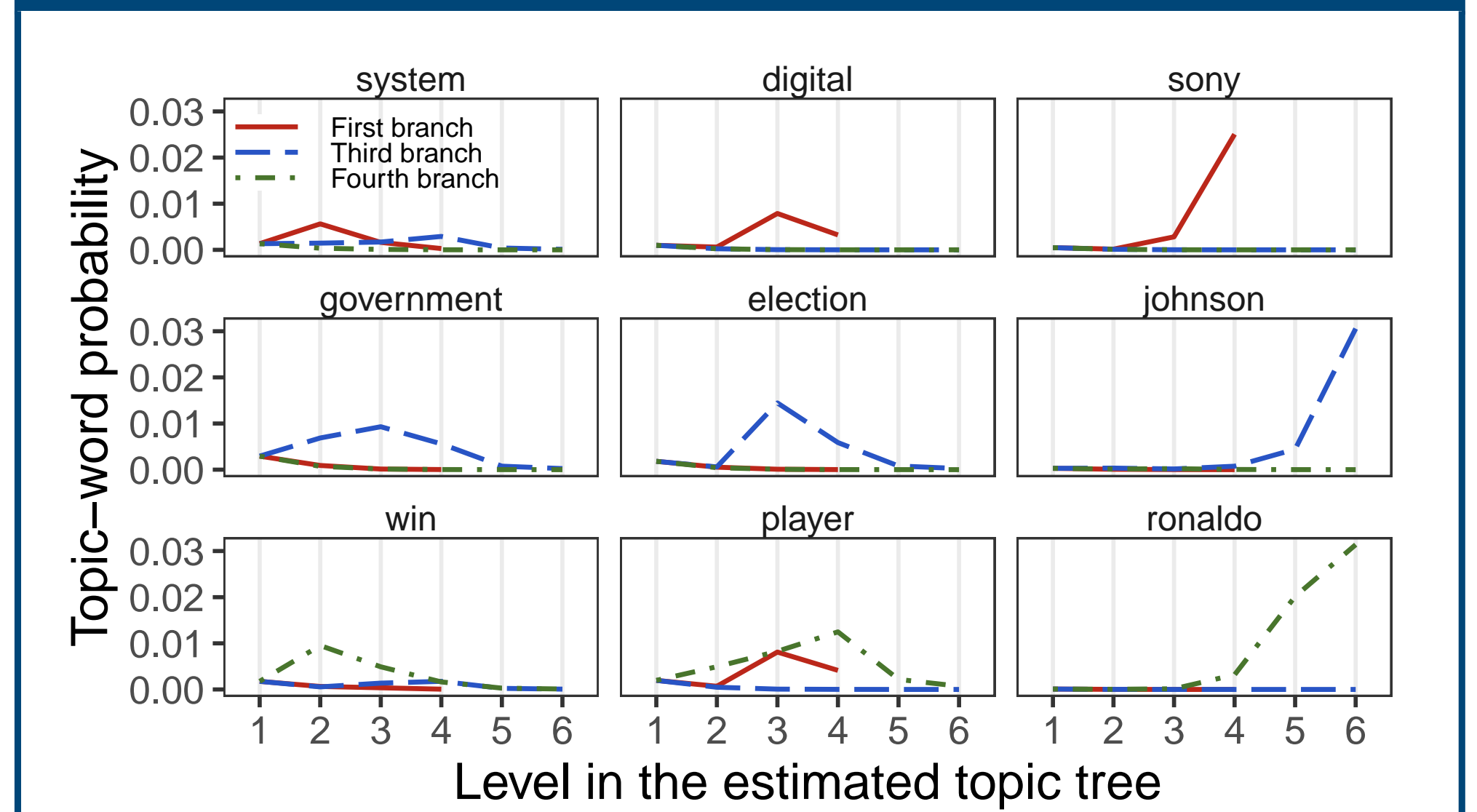
## Experiments

### Data:

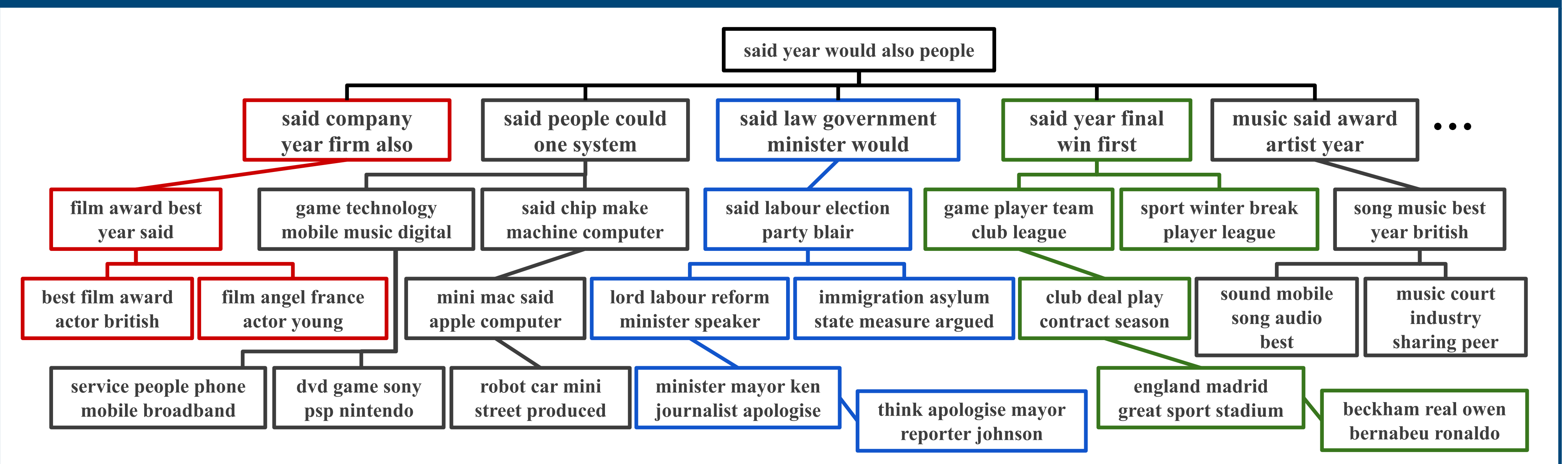
- BBC News: 2,225 documents in five topic areas from the BBC news website
- 20News: a collection of 18,828 posts from 20 USENET newsgroups
- Wikipedia: 50,153 English articles randomly sampled from ten main categories

Setup: Comparing topics with at least 100 assigned words (no truncation in parentheses)

## Topic-Word Probabilities (BBC, L6)



## Top Words of Selected Topics from the Estimated Topic Tree (BBC, L6)



## Evaluation Metrics

Model	Max Lvl.	Tree Diversity ( $\uparrow$ )			Topic Uniqueness ( $\uparrow$ )			Average Overlap ( $\downarrow$ )			# of Topics		
		BBC	20News	Wiki	BBC	20News	Wiki	BBC	20News	Wiki	BBC	20News	Wiki
ihLDA	3	2.24	<b>2.88</b>	<b>2.63</b>	<b>0.60</b>	<b>0.82</b>	<b>0.66</b>	0.28	0.11	0.16	38	27	17
	$\geq 4$	(2.24)	(2.86)	(2.49)	(0.60)	(0.80)	(0.63)	(0.28)	(0.14)	(0.19)	(38)	(31)	(18)
nCRP rCRP	3	<b>2.53</b>	<b>2.88</b>	2.50	0.55	0.76	0.65	0.26	0.12	0.15	85	67	73
		(2.54)	(2.80)	(2.51)	(0.49)	(0.51)	(0.63)	(0.30)	(0.38)	(0.16)	(134)	(203)	(101)
TSNTM	3	1.92	2.16	–	0.36	0.32	–	<b>0.03</b>	0.02	–	517	2108	–
nTSNTM	3	0.15	–	–	0.01	–	–	0.53	–	–	278	–	–
TSNTM	3	1.98	2.54	2.47	0.43	0.80	0.64	0.26	0.09	0.06	22	41	44
nTSNTM	3	2.11	2.57	2.34	0.46	0.68	0.60	0.09	<b>0.01</b>	<b>0.02</b>	68	81	111

- Tree Diversity: uniqueness of topics while considering the importance of the parent topics
- Topic Uniqueness [3]: uniqueness of all topics
- Average Overlap [3]: the average repetition rate of the top words (less overlap does not necessarily mean better interpretability)
- The results without truncation are shown in parentheses for the proposed model.